# Introduction to Machine Learning
## Measures of Accuracy

Andres Mendez-Vazquez

August 21, 2020

# Outline

# Outline

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

## Going back to the probability models

We might think in the machine to be learned as a function $g(x|D)$.....

- Something as curve fitting...

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

## Going back to the probability models

We might think in the machine to be learned as a function $g(\boldsymbol{x}|\mathcal{D})$....

- Something as curve fitting...

Under a data set

$$\mathcal{D} = \{(\boldsymbol{x}_i, y_i)|i = 1, 2, \ldots, N\} \tag{1}$$

Remark: Where the $\boldsymbol{x}_i \sim p(\boldsymbol{x}|\Theta)$!!!

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

## Going back to the probability models

We might think in the machine to be learned as a function $g(\boldsymbol{x}|\mathcal{D})$....

- Something as curve fitting...

Under a data set

$$\mathcal{D} = \{(\boldsymbol{x}_i, y_i)|i = 1, 2, \ldots, N\} \tag{1}$$

Remark: Where the $\boldsymbol{x}_i \sim p(\boldsymbol{x}|\Theta)$!!!

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

## Going back to the probability models

We might think in the machine to be learned as a function $g\left(\boldsymbol{x}|\mathcal{D}\right)$....

- Something as curve fitting...

## Under a data set

$$\mathcal{D} = \{(\boldsymbol{x}_i, y_i) \,|\, i = 1, 2, ..., N\} \tag{1}$$

Remark: Where the $x_i \sim p\left(x|\Theta\right)$!!!

# Introduction

## What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

## Going back to the probability models

We might think in the machine to be learned as a function $g\left(\boldsymbol{x}|\mathcal{D}\right)$....

- Something as curve fitting...

## Under a data set

$$\mathcal{D} = \{(\boldsymbol{x}_i, y_i) \,|\, i = 1, 2, ..., N\} \tag{1}$$

Remark: Where the $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$!!!

# Thus, we have that

## Two main functions

- A function $g(\boldsymbol{x}|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|\boldsymbol{x}]$ the optimal regression...

# Thus, we have that

## Two main functions

- A function $g\left(x|\mathcal{D}\right)$ obtained using some algorithm!!!
- $E\left[y|x\right]$ the optimal regression...

Important

The key factor here is the dependence of the approximation on $\mathcal{D}$

# Thus, we have that

## Two main functions

- A function $g(\boldsymbol{x}|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|\boldsymbol{x}]$ the optimal regression...

## Important

The key factor here is the dependence of the approximation on $\mathcal{D}$.

## Why?

The approximation may be very good for a specific training data set but very bad for another

- This is the reason of studying fusion of information at decision level...

# Thus, we have that

## Two main functions

- A function $g\left(\boldsymbol{x}|\mathcal{D}\right)$ obtained using some algorithm!!!
- $E\left[y|\boldsymbol{x}\right]$ the optimal regression...

## Important

The key factor here is the dependence of the approximation on $\mathcal{D}$.

## Why?

The approximation may be very good for a specific training data set but very bad for another.

- This is the reason of studying fusion of information at decision level...

# Thus, we have that

## Two main functions
- A function $g(\boldsymbol{x}|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|\boldsymbol{x}]$ the optimal regression...

## Important
The key factor here is the dependence of the approximation on $\mathcal{D}$.

## Why?
The approximation may be very good for a specific training data set but very bad for another.
- This is the reason of studying fusion of information at decision level...

# Outline

# How do we measure the difference

## We have that

$$Var(X) = E((X - \mu)^2)$$

# How do we measure the difference

## We have that

$$Var(X) = E((X - \mu)^2)$$

## We can do that for our data

$$Var_{\mathcal{D}}\left(g\left(\boldsymbol{x}|\mathcal{D}\right)\right) = E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

# How do we measure the difference

**We have that**

$$Var(X) = E((X - \mu)^2)$$

**We can do that for our data**

$$Var_{\mathcal{D}}(g(\boldsymbol{x}|\mathcal{D})) = E_D\left((g(\boldsymbol{x}|\mathcal{D}) - E[y|\boldsymbol{x}])^2\right)$$

**Now, if we add and subtract**

$$E_D[g(\boldsymbol{x}|\mathcal{D})] \tag{2}$$

Remark: The expected output of the machine $y(x|\mathcal{D})$

# How do we measure the difference

## We have that

$$Var(X) = E((X - \mu)^2)$$

## We can do that for our data

$$Var_{\mathcal{D}}\left(g\left(\boldsymbol{x}|\mathcal{D}\right)\right) = E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

## Now, if we add and subtract

$$E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] \tag{2}$$

Remark: The expected output of the machine $g\left(\boldsymbol{x}|\mathcal{D}\right)$

# Thus, we have that

## Or Original variance

$$Var_{\mathcal{D}}\left(g\left(\boldsymbol{x}|\mathcal{D}\right)\right) = E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

# Thus, we have that

## Or Original variance

$$Var_{\mathcal{D}}\left(g\left(\boldsymbol{x}|\mathcal{D}\right)\right) = E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$
$$= E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] + E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$
$$= E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2 + ...\right.$$
$$...2\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right) + ...$$
$$...\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

## Finally

$$E_D\left(\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)\right)\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)\right) = ? \qquad (3)$$

# Thus, we have that

## Or Original variance

$$Var_{\mathcal{D}}\left(g\left(\boldsymbol{x}|\mathcal{D}\right)\right) = E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

$$= E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] + E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

$$= E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2 + ...\right.$$

$$...2\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)\right)\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right) + ...$$

$$...\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2\right)$$

## Finally

$$E_D\left(\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)\right)\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)\right) = ? \qquad (3)$$

# Thus, we have that

## Or Original variance

$$Var_{\mathcal{D}} \left( g \left( \boldsymbol{x} | \mathcal{D} \right) \right) = E_D \left( \left( g \left( \boldsymbol{x} | \mathcal{D} \right) - E \left[ y | \boldsymbol{x} \right] \right)^2 \right)$$

$$= E_D \left( \left( g \left( \boldsymbol{x} | \mathcal{D} \right) - E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] + E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] - E \left[ y | \boldsymbol{x} \right] \right)^2 \right)$$

$$= E_D \left( \left( g \left( \boldsymbol{x} | \mathcal{D} \right) - E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] \right)^2 + ... \right.$$

$$...2 \left( \left( g \left( \boldsymbol{x} | \mathcal{D} \right) - E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] \right) \right) \left( E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] - E \left[ y | \boldsymbol{x} \right] \right) + ...$$

$$... \left( E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] - E \left[ y | \boldsymbol{x} \right] \right)^2 \right)$$

## Finally

$$E_D \left( \left( \left( g \left( \boldsymbol{x} | \mathcal{D} \right) - E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] \right) \right) \left( E_D \left[ g \left( \boldsymbol{x} | \mathcal{D} \right) \right] - E \left[ y | \boldsymbol{x} \right] \right) \right) = ? \qquad (3)$$

# Outline

# We have the Bias-Variance

## Our Final Equation

$$E_D\left((g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right])^2\right) = \underbrace{E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2\right)}_{VARIANCE} + \underbrace{\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2}_{BIAS}$$

# We have the Bias-Variance

## Our Final Equation

$$E_D \left( \left( g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right] \right)^2 \right) = \underbrace{E_D \left( \left( g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] \right)^2 \right)}_{VARIANCE} + \underbrace{\left( E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right] \right)^2}_{BIAS}$$

## Where the variance

It represents the measure of the error between our machine $g\left(\boldsymbol{x}|\mathcal{D}\right)$ and the expected output of the machine under $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$.

# We have the Bias-Variance

## Our Final Equation

$$E_D\left((g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right])^2\right) = \underbrace{E_D\left((g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right])^2\right)}_{VARIANCE} + \underbrace{(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right])^2}_{BIAS}$$

## Where the variance

It represents the measure of the error between our machine $g\left(\boldsymbol{x}|\mathcal{D}\right)$ and the expected output of the machine under $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$.

## Where the bias

It represents the quadratic error between the expected output of the machine under $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$ and the expected output of the optimal regression.

# We have the Bias-Variance

## Our Final Equation

$$E_D\left((g\left(\boldsymbol{x}|\mathcal{D}\right) - E\left[y|\boldsymbol{x}\right])^2\right) = \underbrace{E_D\left(\left(g\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2\right)}_{VARIANCE} + \underbrace{\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2}_{BIAS}$$

## Where the variance

It represents the measure of the error between our machine $g\left(\boldsymbol{x}|\mathcal{D}\right)$ and the expected output of the machine under $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$.

## Where the bias

It represents the quadratic error between the expected output of the machine under $\boldsymbol{x}_i \sim p\left(\boldsymbol{x}|\Theta\right)$ and the expected output of the optimal regression.

# Remarks

## We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

# Remarks

## We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

## The situation is more dire in a finite set of data $\mathcal{D}$

We have then a trade-off:

1. Increasing the bias decreases the variance and vice versa.
2. This is known as the **bias–variance dilemma**.

# Remarks

## We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

## The situation is more dire in a finite set of data $\mathcal{D}$

We have then a trade-off:

1. Increasing the bias decreases the variance and vice versa.

2. This is known as the bias–variance dilemma.

# Remarks

## We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

## The situation is more dire in a finite set of data $\mathcal{D}$

We have then a trade-off:

1. Increasing the bias decreases the variance and vice versa.
2. This is known as the **bias–variance dilemma**.

# Similar to...

## Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number $N$, the model will fit the idiosyncrasies of the specific data set.

# Similar to...

## Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number $N$, the model will fit the idiosyncrasies of the specific data set.

## Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

# Similar to...

## Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number $N$, the model will fit the idiosyncrasies of the specific data set.

## Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

## Furthermore

If $N$ grows we can have a more complex model to be fitted which reduces bias and ensures low variance.

- However, $N$ is always finite!!!

# Similar to...

## Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number $N$, the model will fit the idiosyncrasies of the specific data set.

## Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

## Furthermore

If $N$ grows we can have a more complex model to be fitted which reduces bias and ensures low variance.

- However, $N$ is always finite!!!

# Thus

## You always need to compromise

However, you always have some a priori knowledge about the data

## Allowing you to impose restrictions

Lowering the bias and the variance

## Nevertheless

We have the following example to grasp better the bothersome bias–variance dilemma.

# Thus

## You always need to compromise
However, you always have some a priori knowledge about the data

## Allowing you to impose restrictions
Lowering the bias and the variance

## Nevertheless
We have the following example to grasp better the bothersome bias–variance dilemma.

# Thus

## You always need to compromise
However, you always have some a priori knowledge about the data

## Allowing you to impose restrictions
Lowering the bias and the variance

## Nevertheless
We have the following example to grasp better the bothersome **bias–variance dilemma**.

# For this

We know that

The optimum regressor is $E[y|x] = f(x)$

Furthermore

Assume that the randomness in the different training sets, $\mathcal{D}$, is due to the $y_i$'s (Affected by noise), while the respective points, $x_i$, are fixed.

# For this

## Assume

The data is generated by the following function

$$y = f(x) + \epsilon,$$
$$\epsilon \sim \mathcal{N}\left(0, \sigma_\epsilon^2\right)$$

## We know that

The optimum regressor is $E[y|x] = f(x)$

## Furthermore

Assume that the randomness in the different training sets, $\mathcal{D}$, is due to the $y_i$'s (Affected by noise), while the respective points, $x_i$, are fixed.

# For this

**Assume**

The data is generated by the following function
$$y = f(x) + \epsilon,$$
$$\epsilon \sim \mathcal{N}\left(0, \sigma_\epsilon^2\right)$$

**We know that**

The optimum regressor is $E[y|x] = f(x)$

**Furthermore**

Assume that the randomness in the different training sets, $\mathcal{D}$, is due to the $y_i$'s (Affected by noise), while the respective points, $x_i$, are fixed.

# Outline

# Sampling the Space

# Case 1

Choose the estimate of $f(x)$, $g(x|\mathcal{D})$, to be independent of $\mathcal{D}$

For example, $g(x) = w_1 x + w_0$

For example, the points are spread around $(x, f(x))$

# Case 1

For example, $g(x) = w_1 x + w_0$

**For example, the points are spread around $(x, f(x))$**

# Case 1

## Since $g(x)$ is fixed

$$E_{\mathcal{D}}\left[g\left(x|\mathcal{D}\right)\right] = g\left(x|\mathcal{D}\right) \equiv g\left(x\right) \tag{4}$$

## With

$$Var_{\mathcal{D}}\left[g\left(x|\mathcal{D}\right)\right] = 0 \tag{5}$$

## On the other hand

Because $g(x)$ was chosen arbitrarily the expected bias must be large

$$\underbrace{\left(E_{\mathcal{D}}\left[g\left(x|\mathcal{D}\right)\right] - E\left[y|x\right]\right)^2}_{BIAS} \tag{6}$$

# Case 1

## Since $g(x)$ is fixed

$$E_{\mathcal{D}}[g(x|\mathcal{D})] = g(x|\mathcal{D}) \equiv g(x) \tag{4}$$

## With

$$Var_{\mathcal{D}}[g(x|\mathcal{D})] = 0 \tag{5}$$

## On the other hand

Because $g(x)$ was chosen arbitrarily the expected bias must be large

$$\underbrace{(E_{\mathcal{D}}[g(x|\mathcal{D})] - E[y|x])^2}_{BIAS} \tag{6}$$

# Case 1

**Since $g(x)$ is fixed**

$$E_{\mathcal{D}}\left[g\left(x|\mathcal{D}\right)\right] = g\left(x|\mathcal{D}\right) \equiv g\left(x\right) \tag{4}$$

**With**

$$Var_{\mathcal{D}}\left[g\left(x|\mathcal{D}\right)\right] = 0 \tag{5}$$

**On the other hand**

Because $g\left(x\right)$ was chosen arbitrarily the expected bias must be large.

$$\underbrace{\left(E_D\left[g\left(\boldsymbol{x}|\mathcal{D}\right)\right] - E\left[y|\boldsymbol{x}\right]\right)^2}_{BIAS} \tag{6}$$

# Case 2

## In the other hand

Now, $g_1(x)$ corresponds to a polynomial of high degree so it can pass through each training point in $\mathcal{D}$.

Example of $g_1(x)$

# Case 2

## In the other hand

Now, $g_1(x)$ corresponds to a polynomial of high degree so it can pass through each training point in $\mathcal{D}$.

## Example of $g_1(x)$

# Case 2

## Due to the zero mean of the noise source

$$E_D\left[g_1\left(\boldsymbol{x}|\mathcal{D}\right)\right] = f\left(x\right) = E\left[y|x\right] \text{ for any } x = x_i \tag{7}$$

Remark: At the training points the bias is zero.

However the variance increases

$$E_D\left[\left(y_1\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[y_1\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2\right] = E_D\left[\left(f\left(x\right) + \epsilon - f\left(x\right)\right)^2\right]$$

$$= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, \ldots, N$$

In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source

# Case 2

$$E_D\left[g_1\left(\boldsymbol{x}|\mathcal{D}\right)\right] = f\left(x\right) = E\left[y|x\right] \text{ for any } x = x_i \tag{7}$$

Remark: At the training points the bias is zero.

However the variance increases

$$E_D\left[\left(g_1\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g_1\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2\right] = E_D\left[\left(f\left(x\right) + \epsilon - f\left(x\right)\right)^2\right]$$
$$= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, ..., N$$

In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source

# Case 2

$$E_D\left[g_1\left(\boldsymbol{x}|\mathcal{D}\right)\right] = f\left(x\right) = E\left[y|x\right] \text{ for any } x = x_i \tag{7}$$

Remark: At the training points the bias is zero.

## However the variance increases

$$E_D\left[\left(g_1\left(\boldsymbol{x}|\mathcal{D}\right) - E_D\left[g_1\left(\boldsymbol{x}|\mathcal{D}\right)\right]\right)^2\right] = E_D\left[\left(f\left(x\right) + \epsilon - f\left(x\right)\right)^2\right]$$
$$= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, ..., N$$

## In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source.

## First

Everything that has been said so far applies to both the regression and the classification tasks.

## However

Mean squared error is not the best way to measure the power of a classifier.

## Think about

A classifier that sends everything far away of the hyperplane!!! Away from the values $+ - 1$!!!

# Observations

## First

Everything that has been said so far applies to both the regression and the classification tasks.

## However

Mean squared error is not the best way to measure the power of a classifier.

## Think about

A classifier that sends everything far away of the hyperplane!!! Away from the values $+ - 1$!!!

# Observations

## First

Everything that has been said so far applies to both the regression and the classification tasks.

## However

Mean squared error is not the best way to measure the power of a classifier.

## Think about

A classifier that sends everything far away of the hyperplane!!! Away from the values $+-1$!!!

# Outline

# Sooner of Latter you need to know how efficient is your algorithm

## Thus, we need a measures of accuracy

Thus, we begin with the classic classifier for two classes

# Sooner of Latter you need to know how efficient is your algorithm

## Thus, we need a measures of accuracy

Thus, we begin with the classic classifier for two classes



## Here

A dataset used for performance evaluation is called a **test dataset.**

# Therefore

> **It is a good idea to build a measure of performance**
>
> For this, we can use the idea of error in statistics.

# Therefore

## It is a good idea to build a measure of performance

For this, we can use the idea of error in statistics.

## Do you remember?

# Outline

# $\alpha$ error

## Definition (Type I Error - False Positive)

$\alpha$ is the probability that the test will lead to the rejection of the hypothesis $H_0$ when that hypothesis is true.

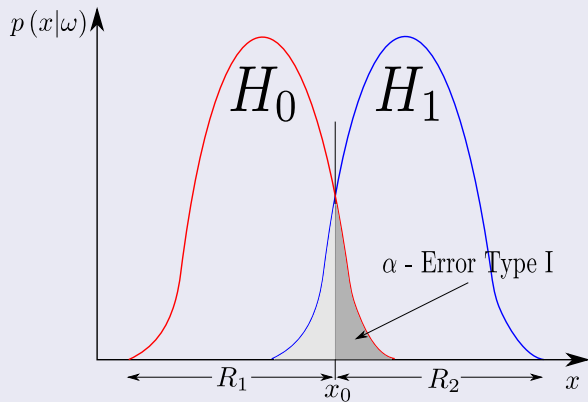# $\alpha$ error

## Definition (Type I Error - False Positive)

$\alpha$ is the probability that the test will lead to the rejection of the hypothesis $H_0$ when that hypothesis is true.

## Example

1. $H_0$ : "You have a device that produce circuits with no error"

2. You have a device that fails $\alpha = 0.05$ meaning that it fails 5 of the time.

3. This says that you ha low chance of a wrong circuit.

# $\alpha$ error

## Definition (Type I Error - False Positive)

$\alpha$ is the probability that the test will lead to the rejection of the hypothesis $H_0$ when that hypothesis is true.

## Example

1. $H_0$ : "You have a device that produce circuits with no error"
2. You have a device that fails $\alpha = 0.05$ meaning that it fails $5$ of the time.
3. This says that you ha low chance of a wrong circuit.

# $\alpha$ error

## Definition (Type I Error - False Positive)

$\alpha$ is the probability that the test will lead to the rejection of the hypothesis $H_0$ when that hypothesis is true.

## Example

1. $H_0$ : "You have a device that produce circuits with no error"
2. You have a device that fails $\alpha = 0.05$ meaning that it fails $5$ of the time.
3. This says that you ha low chance of a wrong circuit.

# Basically

# $\beta$ error

## Definition (Type II Error - False Negative)

$\beta$ is the probability that the test will lead to the rejection of the hypothesis $H_1$ when that hypothesis is true.

# $\beta$ error

### Definition (Type II Error - False Negative)

$\beta$ is the probability that the test will lead to the rejection of the hypothesis $H_1$ when that hypothesis is true.

### Example

1. $H_1$: "Adding fluoride to toothpaste protects against cavities."

2. Then $\beta = 0.05$ meaning that you have a chance of 5 of the time.

3. This says that you ha low chance of having a cavity using fluoride in the water

# $\beta$ error

## Definition (Type II Error - False Negative)

$\beta$ is the probability that the test will lead to the rejection of the hypothesis $H_1$ when that hypothesis is true.

## Example

1. $H_1$: "Adding fluoride to toothpaste protects against cavities."
2. Then $\beta = 0.05$ meaning that you have a chance of $5$ of the time.
3. This says that you ha low chance of having a cavity using fluoride in the water

# $\beta$ error

## Definition (Type II Error - False Negative)

$\beta$ is the probability that the test will lead to the rejection of the hypothesis $H_1$ when that hypothesis is true.

## Example

1. $H_1$: "Adding fluoride to toothpaste protects against cavities."
2. Then $\beta = 0.05$ meaning that you have a chance of $5$ of the time.
3. This says that you ha low chance of having a cavity using fluoride in the water.

# Outline

# This can be seen as a table

## Confusion Matrix

| Table of error types | | Null Hypothesis $H_0$ | |
|---|---|---|---|
| | | True | False |
| Decision about $H_0$ | Reject | Type I Error - $\alpha$ **False Positive** | Correct Inference **True Positive** |
| | Fail to reject | Correct Inference **True Negative** | Type II Error - $\beta$ **False Negative** |

# In the case of two classes, we have

## We have the following

|  |  | Actual Class | |
|---|---|---|---|
|  |  | Positive | Negative |
| Predicted | Positive | **True Positive (TP)** | **False Positives (FP)** |
| Classes | Negative | **False Negatives (FN)** | **True Negatives (TN)** |

# Outline

# Accuracy

## Definition

The proportion of getting correct classification of the Positive and Negative classes.

Thus

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

Problem - accuracy assumes equal cost for both kinds of errors

Is 99% accuracy good, bad or terrible? It depends on the problem.

# Accuracy

### Definition

The proportion of getting correct classification of the Positive and Negative classes.

### Thus

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

Problem: accuracy assumes equal cost for both kinds of errors

Is 99% accuracy good, bad or terrible? It depends on the problem.

# Accuracy

**Definition**

The proportion of getting correct classification of the Positive and Negative classes.

**Thus**

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

**Problem - accuracy assumes equal cost for both kinds of errors**

Is 99% accuracy good, bad or terrible? It depends on the problem.

# True Positive Rate

### Definition

True Positive Rate is the proportion of getting a correct classification of the Positive Class vs the True Positive and False Negatives.

$$\text{True Positive Rate} = \frac{TP}{TP + FN}$$

# True Positive Rate

## Also called
Sensitivity or Recall Rate

## Defined as
True Positive Rate is the proportion of getting a correct classification of the Positive Class vs the True Positive and False Negatives.

$$\text{True Positive Rate} = \frac{TP}{TP + FN}$$

# True Negative Rate

## Also known as
Specificity

## Defined as
It is the proportion of True Negative vs the elements classified as True negatives.

$$\text{True Negative Rate} = \frac{TN}{FP + TN}$$

# True Negative Rate

## Also known as

Specificity

## Defined as

It is the proportion of True Negative vs the elements classified as True negatives.

$$\text{True Negative Rate} = \frac{TN}{FP + TN}$$

# Precision

## Also known as

Positive Predictive Value

## Defined as

The proportion of the elements classified as true positive vs the total of all the real true positives.

$$\text{Precision Predicted Value} = \frac{TP}{FP + TP}$$

# Precision

## Also known as

Positive Predictive Value

## Defined as

The proportion of the elements classified as true positive vs the total of all the real true positives.

$$\text{Precision Predicted Value} = \frac{TP}{FP + TP}$$

# Significance Level

## Also known as

False Positive Rate.

Defined as

False Positive Rate is the probability of getting an incorrect classification of the Positive Class vs the True Negative and the False Positive.

$$\text{False positive rate} = \frac{FP}{TN + FP}$$

# Significance Level

## Also known as

False Positive Rate.

## Defined as

False Positive Rate is the probability of getting an incorrect classification of the Positive Class vs the True Negative and the False Positive.

$$\text{False positive rate} = \frac{FP}{TN + FP}$$

# Outline

# We can do better than these simple measures of accuracy

## Given these initial measures of validity

it is possible to obtain a more precise model evaluation, the ROC curves.

## The ROC Curves plot

It is a model-wide evaluation measure that is based on two basic
evaluation measures:

1. **Specificity** is a performance measure of the whole negative part of a
   dataset.

2. **Sensitivity** is a performance measure of the whole positive part.

# We can do better than these simple measures of accuracy

## Given these initial measures of validity

it is possible to obtain a more precise model evaluation, the ROC curves.

## The ROC Curves plot

It is a model-wide evaluation measure that is based on two basic evaluation measures:

1. **Specificity** is a performance measure of the whole negative part of a dataset.
2. **Sensitivity** is a performance measure of the whole positive part.

# What the ROC Curves uses

## We have a plot where

The ROC plot uses specificity on the $x$-axis and sensitivity on the $y$-axis.

# What the ROC Curves uses

## We have a plot where

The ROC plot uses specificity on the $x$-axis and sensitivity on the $y$-axis.

## Basically

False Positive Rate (FPR) is identical with specificity, and True Positive Rate (TPR) is identical with sensitivity.

# What the ROC Curves uses

## We have a plot where

The ROC plot uses specificity on the $x$-axis and sensitivity on the $y$-axis.

## Basically

False Positive Rate (FPR) is identical with specificity, and True Positive Rate (TPR) is identical with sensitivity.

## Then

1. A ROC curve is created by connecting all ROC points of a classier in the ROC space.
2. Two adjacent ROC points can be connected by a straight line.
3. The curve starts at (0.0, 0.0) and ends at (1.0, 1.0).

# What the ROC Curves uses

## We have a plot where

The ROC plot uses specificity on the $x$-axis and sensitivity on the $y$-axis.

## Basically

False Positive Rate (FPR) is identical with specificity, and True Positive Rate (TPR) is identical with sensitivity.

## Then

1. A ROC curve is created by connecting all ROC points of a classier in the ROC space.
2. Two adjacent ROC points can be connected by a straight line.
3. The curve starts at (0.0, 0.0) and ends at (1.0, 1.0)

# What the ROC Curves uses

## We have a plot where

The ROC plot uses specificity on the $x$-axis and sensitivity on the $y$-axis.

## Basically

False Positive Rate (FPR) is identical with specificity, and True Positive Rate (TPR) is identical with sensitivity.

## Then

1. A ROC curve is created by connecting all ROC points of a classier in the ROC space.
2. Two adjacent ROC points can be connected by a straight line.
3. The curve starts at (0.0, 0.0) and ends at (1.0, 1.0).

# Outline

# For Example

# Outline

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate
that example $i$ is positive; P and N, the number of positive and negative
examples.

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ sorted decreasing by $f$ scores
2. $FP \leftarrow TP = 0; R \leftarrow \{\}; f_{prev} \leftarrow -\infty; i \leftarrow 1$
3. while $i \leq |L_{sorted}|$
4. if $f(i) \neq f_{prev}$ then
5. $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6. $f_{prev} \leftarrow f(i)$
7. if $L_{sorted}$ is a positive example then $TP = TP + 1$
8. else $FP = FP + 1$
9. $i \leftarrow i + 1$
10. $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$, this is $(1,1)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

① $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**

② $FP \leftarrow TP \leftarrow 0; R \leftarrow \{\}; f_{prev} \leftarrow -\infty; i \leftarrow 1$

③ while $i \leq |L_{sorted}|$

④ if $f(i) \neq f_{prev}$ then

⑤ $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$

⑥ $f_{prev} \leftarrow f(i)$

⑦ if $L_{sorted}$ is a positive example then $TP = TP + 1$

⑧ else $FP = FP + 1$

⑨ $i \leftarrow i + 1$

⑩ $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$, this is $(1,1)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle ; fprev \leftarrow -\infty; i \leftarrow 1$
3. while $i \leq |L_{sorted}|$
4.      if $f(i) \neq f_{prev}$ then
5.          $R.append \left( \frac{FP}{N}, \frac{TP}{P} \right)$
6.          $fprev \leftarrow f(i)$
7.      if $L_{sorted}$ is a positive example then $TP = TP + 1$
8.      else $FP = FP + 1$
9.      $i \leftarrow i + 1$
10. $R.append \left( \frac{FP}{N}, \frac{TP}{P} \right)$, this is $(1, 1)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle\rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.       **if** $f(i) \neq f_{prev}$ **then**
5.            $R.append\left( \frac{FP}{N}, \frac{TP}{P} \right)$
6.            $fprev \leftarrow f(i)$
7.            **if** $L_{sorted}$ **is a positive example then** $TP \leftarrow TP + 1$
8.            **else** $FP \leftarrow FP + 1$
9.            $i \leftarrow i + 1$
10.       $R.append\left( \frac{FP}{N}, \frac{TP}{P} \right)$, this is $(1, 1)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$
7.          if $L_{sorted}$ is a positive example then $TP = TP + 1$
8.          else $FP = FP + 1$
9.          $i \leftarrow i + 1$
10.     $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$, this is $(1,1)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$
7.      **if** $L_{sorted}$ **is a positive example then** $TP = TP + 1$

# We have

## Algorithm ROC point generation

> Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

> Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by** $f$ **scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle ; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$
7.      **if** $L_{sorted}$ **is a positive example then** $TP = TP + 1$
8.      **else** $FP = FP + 1$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$
7.      **if** $L_{sorted}$ **is a positive example then** $TP = TP + 1$
8.      **else** $FP = FP + 1$
9.      $i \leftarrow i + 1$

# We have

## Algorithm ROC point generation

Input: $L$, the set of test examples; $f(i)$, the probabilistic classifier estimate that example $i$ is positive; P and N, the number of positive and negative examples.

Ouput: $R$, a list of ROC points increasing by false positive rate.

1. $L_{sorted} \leftarrow L$ **sorted decreasing by $f$ scores**
2. $FP \leftarrow TP \leftarrow 0; R \leftarrow \langle \rangle; fprev \leftarrow -\infty; i \leftarrow 1$
3. **while** $i \leq |L_{sorted}|$
4.      **if** $f(i) \neq f_{prev}$ **then**
5.          $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$
6.          $fprev \leftarrow f(i)$
7.      **if** $L_{sorted}$ **is a positive example then** $TP = TP + 1$
8.      **else** $FP = FP + 1$
9.      $i \leftarrow i + 1$
10. $R.append\left(\frac{FP}{N}, \frac{TP}{P}\right)$, **this is** $(1, 1)$

# Thus

Thus, after generating the ROC Curve it is possible to use several metrics to validate using the ROC curves.

A Partial List is
1. Area Under the Curve (AUC)
2. Equal Error Rate (EER)
3. Likelihood Ratio

# Thus

Thus, after generating the ROC Curve it is possible to use several metrics to validate using the ROC curves.

## A Partial List is

1. Area Under the Curve (AUC)
2. Equal Error Rate (EER)
3. Likelihood Ratio

# Outline

# A Simple Defintion

**We have**

$$AUC = \int ROC\,(p)\,dp = \sum_{i=1}^{N} ROC\left(f\left(\frac{1}{i}\right)\right)\left[\frac{i}{N} - \frac{i-1}{N}\right]$$

This equation has the following meaning

- The probability that a randomly selected observation $X$ from the **positive class** would have a higher score than a randomly selected observation $Y$ from the **negative class**.

$$P\,(X > Y)$$

Thus

The AUC gives the mean **true positive** rate averaged uniformly across the **false positive** rate.

# A Simple Defintion

**We have**

$$AUC = \int ROC\left(p\right) dp = \sum_{i=1}^{N} ROC\left(f\left(\frac{1}{i}\right)\right)\left[\frac{i}{N} - \frac{i-1}{N}\right]$$

**This equation has the following meaning**

- The probability that a randomly selected observation $X$ from the **positive class** would have a higher score than a randomly selected observation $Y$ from the **negative class.**

$$P\left(X > Y\right)$$

**Thus**

The AUC gives the mean **true positive** rate averaged uniformly across the **false positive** rate.

# A Simple Defintion

**We have**

$$AUC = \int ROC\left(p\right) dp = \sum_{i=1}^{N} ROC\left(f\left(\frac{1}{i}\right)\right)\left[\frac{i}{N} - \frac{i-1}{N}\right]$$

**This equation has the following meaning**

- The probability that a randomly selected observation $X$ from the **positive class** would have a higher score than a randomly selected observation $Y$ from the **negative class.**

$$P\left(X > Y\right)$$

**Thus**

The AUC gives the mean **true positive** rate averaged uniformly across the **false positive** rate.

# Outline

# Also known as $F_1$ score

## It is a measure of a test's accuracy

It considers both the precision $P$ and the recall $R$ of the test to compute the score.

## An interesting fact

It computes some average of the information retrieval precision and recall.

# Also known as $F_1$ score

## It is a measure of a test's accuracy

It considers both the precision $P$ and the recall $R$ of the test to compute the score.

## An interesting fact

It computes some average of the information retrieval precision and recall.

# Comparison of Measures

## Something Notable

$$Average = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$Harmonic = \frac{N}{\sum_{i=1}^{N} \frac{1}{x_i}}$$

When $x_i = Precision$ and $x_2 = Recall$

$$Average = \frac{1}{2}(P + R)$$

$$Harmonic = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P + R}$$

# Comparison of Measures

$$Average = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$Harmonic = \frac{N}{\sum_{i=1}^{N} \frac{1}{x_i}}$$

When $x_1 = Precision$ and $x_2 = Recall$

$$Average = \frac{1}{2} \left( P + R \right)$$

$$Harmonic = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P + R}$$

# Thus

**Important**

- The harmonic mean is more intuitive than the arithmetic mean when computing a mean of ratios.

Example

- Suppose that you have a finger print recognition system and its precision and recall be 1.0 and 0.2

# Thus

## Important

- The harmonic mean is more intuitive than the arithmetic mean when computing a mean of ratios.

## Example

- Suppose that you have a finger print recognition system and its precision and recall be 1.0 and 0.2

# How is this Computed?

**Then for Precision and Recall, we have a general function**

$$F_\beta = \frac{(\beta^2 + 1)\,Precision \times Recall}{\beta^2 Precision + Recall}\ \ (0 \le \beta \le +\infty)$$

Thus, for the basic case, $F_1$

$$F_1 = 2\frac{Precision \times Recall}{Precision + Recall}$$

# How is this Computed?

## Then for Precision and Recall, we have a general function

$$F_\beta = \frac{(\beta^2 + 1)\, Precision \times Recall}{\beta^2 Precision + Recall} \quad (0 \leq \beta \leq +\infty)$$

## Thus, for the basic case $F_1$

$$F_1 = 2\frac{Precision \times Recall}{Precision + Recall}$$

# Outline

# What we want

## We want to measure

A quality measure to measure different classifiers (for different parameter values).

We call that as

$$R(f) = E_D[L(y, f(x))].$$ (8)

Example $L(y, f(x)) = \|y - f(x)\|_2^2$

More precisely

For different values $\gamma_j$ of the parameter, we train a classifier $f(x|\gamma_j)$ on the training set.

# What we want

A quality measure to measure different classifiers (for different parameter values).

## We call that as

$$R(f) = E_{\mathcal{D}} \left[ L\left(y, f\left(\boldsymbol{x}\right)\right) \right]. \tag{8}$$

Example: $L\left(y, f\left(x\right)\right) = \|y - f\left(\boldsymbol{x}\right)\|_2^2$

More precisely

For different values $\gamma_y$ of the parameter, we train a classifier $f\left(x|\gamma_y\right)$ on the training set.

# What we want

## We want to measure

A quality measure to measure different classifiers (for different parameter values).

## We call that as

$$R(f) = E_{\mathcal{D}} \left[ L\left(y, f\left(\boldsymbol{x}\right)\right) \right]. \tag{8}$$

Example: $L\left(y, f\left(x\right)\right) = \|y - f\left(\boldsymbol{x}\right)\|_2^2$

## More precisely

For different values $\gamma_j$ of the parameter, we train a classifier $f\left(\boldsymbol{x}|\gamma_j\right)$ on the training set.

# Then, calculate the empirical Risk

**Do you have any ideas?**

Give me your best shot!!!

# Then, calculate the empirical Risk

## Do you have any ideas?

Give me your best shot!!!

## Empirical Risk

We use the validation set to estimate

$$\hat{R}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v} \sum_{i=1}^{N_v} L\left(y_i, f\left(\boldsymbol{x}_i|\gamma\right)\right) \tag{9}$$

# Then, calculate the empirical Risk

## Do you have any ideas?

Give me your best shot!!!

## Empirical Risk

We use the validation set to estimate

$$\hat{R}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v} \sum_{i=1}^{N_v} L\left(y_i, f\left(\boldsymbol{x}_i|\gamma\right)\right) \tag{9}$$

## Thus, you follow the following procedure

1. Select the value $\gamma^*$ which achieves the smallest estimated error.
2. Re-train the classifier with parameter $\gamma^*$ on all data except the test set (i.e. train + validation data)
3. Report error estimate $\hat{R}\left(f\left(x|\gamma\right)\right)$ computed on the test set.

# Then, calculate the empirical Risk

**Empirical Risk**

We use the validation set to estimate

$$\hat{R}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v} \sum_{i=1}^{N_v} L\left(y_i, f\left(\boldsymbol{x}_i|\gamma\right)\right) \tag{9}$$

**Thus, you follow the following procedure**

1. Select the value $\gamma^*$ which achieves the smallest estimated error.
2. Re-train the classifier with parameter $\gamma^*$ on all data except the test set (i.e. train + validation data).
3. Report error estimate $\hat{R}\left(f\left(x|\gamma\right)\right)$ computed on the test set.

# Then, calculate the empirical Risk

## Do you have any ideas?
Give me your best shot!!!

## Empirical Risk
We use the validation set to estimate

$$\hat{R}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v}\sum_{i=1}^{N_v} L\left(y_i, f\left(\boldsymbol{x}_i|\gamma\right)\right) \tag{9}$$

## Thus, you follow the following procedure
1. Select the value $\gamma^*$ which achieves the smallest estimated error.
2. Re-train the classifier with parameter $\gamma^*$ on all data except the test set (i.e. train + validation data).
3. Report error estimate $\hat{R}\left(f\left(x|\gamma_i\right)\right)$ computed on the test set.

## Idea

### Something Notable

- Each of the **error estimates computed on validation set** is computed from a single example of a trained classifier.
  - Can we improve the estimate?

# Idea

## Something Notable

- Each of the **error estimates computed on validation set** is computed from a single example of a trained classifier.
  - Can we improve the estimate?

## $K$-fold Cross Validation

To estimate the risk of a classifier $f$:

1. Split data into $K$ equally sized parts (called "folds"), $N_v$.

2. Train an instance $f_k$ of the classifier, using all folds except fold $k$ as training data.

3. Compute the Cross Validation (CV) estimate:

$$R_{CV}(f(x|\gamma)) = \frac{1}{N_v} \sum_{k=1}^{N_v} L\left(y_i, f_k\left(x_{k(i)}|\gamma\right)\right) \qquad (10)$$

where $k(i)$ is the fold containing $x_i$.

# Idea

## Something Notable

- Each of the **error estimates computed on validation set** is computed from a single example of a trained classifier.
  - Can we improve the estimate?

## $K$-fold Cross Validation

To estimate the risk of a classifier $f$:

1. Split data into $K$ equally sized parts (called "folds"), $N_v$.

2. Train an instance $f_k$ of the classifier, using all folds except fold $k$ as training data.

3. Compute the Cross Validation (CV) estimate:

$$R_{CV}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v} \sum_{i=1}^{N_v} L\left(y_i, f_k\left(x_{k(i)}|\gamma\right)\right) \quad (10)$$

where $k\left(i\right)$ is the fold containing $x_i$.

# Idea

## Something Notable

- Each of the **error estimates computed on validation set** is computed from a single example of a trained classifier.
  - Can we improve the estimate?

## $K$-fold Cross Validation

To estimate the risk of a classifier $f$:

1. Split data into $K$ equally sized parts (called "folds"), $N_v$.
2. Train an instance $f_k$ of the classifier, using all folds except fold $k$ as training data.

3. Compute the Cross Validation (CV) estimate:

$$R_{CV}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v}\sum_{k=1}^{N_v} L\left(y_i, f_k\left(x_{k(i)}|\gamma\right)\right) \tag{10}$$

where $k\left(i\right)$ is the fold containing $x_i$.

# Idea

## Something Notable

- Each of the **error estimates computed on validation set** is computed from a single example of a trained classifier.
  - Can we improve the estimate?

## $K$-fold Cross Validation

To estimate the risk of a classifier $f$:

1. Split data into $K$ equally sized parts (called "folds"), $N_v$.
2. Train an instance $f_k$ of the classifier, using all folds except fold $k$ as training data.
3. Compute the Cross Validation (CV) estimate:

$$\hat{R}_{CV}\left(f\left(x|\gamma\right)\right) = \frac{1}{N_v} \sum_{k=1}^{N_v} L\left(y_i, f_k\left(\boldsymbol{x}_{k(i)}|\gamma\right)\right) \tag{10}$$

where $k\left(i\right)$ is the fold containing $\boldsymbol{x}_i$.

# Example

## $K = 5, k = 3$

| Train | Train | Testing | Train | Train |
|-------|-------|---------|-------|-------|
| 1 | 2 | 3 | 4 | 5 |

# Example

| Train | Train | Testing | Train | Train |
|-------|-------|---------|-------|-------|
| 1     | 2     | 3       | 4     | 5     |

## Actually, we have

Cross validation procedure does not involve the test data.

| SPLIT All Train Set | |
|-------------------------------|------|
| Train Data + Validation Data | Test |

# Outline

# How to choose $K$

## Extremal cases

- $K = N$, called leave one out cross validation (loocv)
- $K = 2$

An often-cited problem with loocv is that we have to train many $(= N)$ classifiers, but there is also a deeper problem.

# How to choose $K$

- $K = N$, called leave one out cross validation (loocv)
- $K = 2$

An often-cited problem with loocv is that we have to train many ($= N$) classifiers, but there is also a deeper problem.

Argument 1: $K$ should be small, e.g. $K = 2$

1. Unless we have a lot of data, variance between two distinct training sets may be considerable.

2. Important concept: By removing substantial parts of the sample in turn and at random, we can simulate this variance.

3. By removing a single point (loocv), we cannot make this variance visible.

# How to choose $K$

## Extremal cases

- $K = N$, called leave one out cross validation (loocv)
- $K = 2$

An often-cited problem with loocv is that we have to train many $(= N)$ classifiers, but there is also a deeper problem.

# How to choose $K$

## Extremal cases

- $K = N$, called leave one out cross validation (loocv)
- $K = 2$

An often-cited problem with loocv is that we have to train many $(= N)$ classifiers, but there is also a deeper problem.

## Argument 1: $K$ should be small, e.g. $K = 2$

1. Unless we have a lot of data, variance between two distinct training sets may be considerable.

2. Important concept: By removing substantial parts of the sample in turn and at random, we can simulate this variance.

3. By removing a single point (loocv), we cannot make this variance visible.

# How to choose $K$

**Argument 1: $K$ should be small, e.g. $K = 2$**

1. Unless we have a lot of data, variance between two distinct training sets may be considerable.
2. Important concept: By removing substantial parts of the sample in turn and at random, we can simulate this variance.
3. By removing a single point (loocv), we cannot make this variance visible.

# How to choose $K$

## Extremal cases

- $K = N$, called leave one out cross validation (loocv)
- $K = 2$

An often-cited problem with loocv is that we have to train many $(= N)$ classifiers, but there is also a deeper problem.

## Argument 1: $K$ should be small, e.g. $K = 2$

1. Unless we have a lot of data, variance between two distinct training sets may be considerable.
2. Important concept: By removing substantial parts of the sample in turn and at random, we can simulate this variance.
3. By removing a single point (loocv), we cannot make this variance visible.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.

2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.

3. This way, we will systematically overestimate the risk.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.
2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.
3. This way, we will systematically overestimate the risk.

## Common recommendation: $K = 5$ to $K = 10$

Intuition:

1. $K = 10$ means number of samples removed from training is one order of magnitude below training sample size.
2. This should not weaken the classifier considerably, but should be large enough to make measure variance effects.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.
2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.
3. This way, we will systematically overestimate the risk.

## Common recommendation: $K = 5$ to $K = 10$

Intuition:

1. $K = 10$ means number of samples removed from training is one order of magnitude below training sample size.
2. This should not weaken the classifier considerably, but should be large enough to make measure variance effects.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.
2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.
3. This way, we will systematically overestimate the risk.

## Common recommendation: $K = 5$ to $K = 10$

Intuition:

1. $K = 10$ means number of samples removed from training is one order of magnitude below training sample size.
2. This should not weaken the classifier considerably, but should be large enough to make measure variance effects.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.
2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.
3. This way, we will systematically overestimate the risk.

## Common recommendation: $K = 5$ to $K = 10$

Intuition:

1. $K = 10$ means number of samples removed from training is one order of magnitude below training sample size.
2. This should not weaken the classifier considerably, but should be large enough to make measure variance effects.

# How to choose $K$

## Argument 2: $K$ should be large, e.g. $K = N$

1. Classifiers generally perform better when trained on larger data sets.
2. A small $K$ means we substantially reduce the amount of training data used to train each $f_k$, so we may end up with weaker classifiers.
3. This way, we will systematically overestimate the risk.

## Common recommendation: $K = 5$ to $K = 10$

Intuition:

1. $K = 10$ means number of samples removed from training is one order of magnitude below training sample size.
2. This should not weaken the classifier considerably, but should be large enough to make measure variance effects.