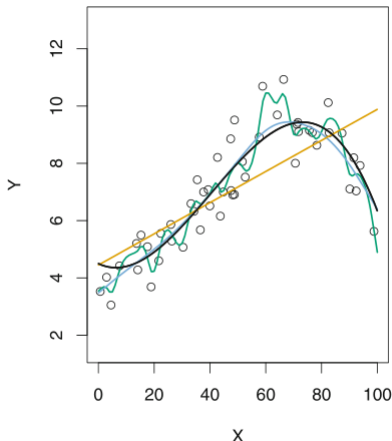


Prediction and Classification

Statistical Learning

Dr. Markus Heinrich

Sample vs Population and True vs Estimated Predictor Function



- What is the optimal prediction given a loss function?
- How is the prediction error decomposed?
- How can I minimize the prediction error?

Prediction and Classification

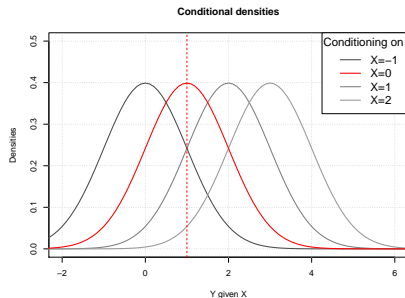
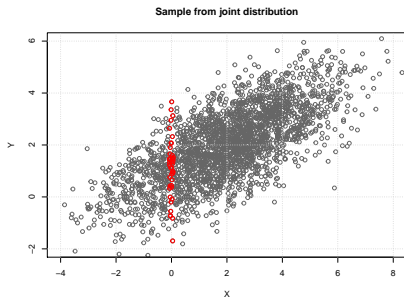
- 1 Prediction errors and quantification
- 2 Classification errors
- 3 Learning for prediction
- 4 Up next

Outline

- 1 Prediction errors and quantification
- 2 Classification errors
- 3 Learning for prediction
- 4 Up next

There will be errors

Conditionally on \mathbf{X} , Y has a distribution and its outcomes (continuous or discrete) are random.



At the same time, \hat{Y} is a deterministic function of \mathbf{X} and therefore prediction (or classification) errors will appear by construction.

To minimize the (impact of the) errors, we need to quantify them.

Let us discuss prediction first

For **continuous** variables,

- The probability of $Y = \hat{Y}$ is zero.
- A good prediction lies as close as possible to Y **on average**.
- Seek a *good predictor function* $f(\mathbf{X}) = \hat{Y}$ for predicting Y given values of the input \mathbf{X} .

But what is good / close?

We shall measure this with the help of a loss function $\mathcal{L}(Y, f(\mathbf{X}), \dots)$ for penalizing errors $u = (Y - \hat{Y})$ in prediction.

Loss function

Definition

A loss function $\mathcal{L}(\cdot)$ is a quasi-convex function minimized at 0.

In other words, \mathcal{L} is increasing for positive arguments and decreasing for negative arguments.

- The most common loss function is the squared error loss (MSE).
- But asymmetric loss functions may also be encountered, say asymmetric linear or asymmetric quadratic.

Typically, we pick the (optimal) prediction that minimizes the expected loss!

$$\hat{Y}^* = \arg \min_y E(\mathcal{L}(Y, y)).$$

What is the optimal prediction?

General approach:

- ① State your loss (objective) function $\mathcal{L}(Y, \hat{Y})$
- ② Minimize expected loss $E[\mathcal{L}(Y, \hat{Y})]$:

- First derivative:

$$\frac{d}{d\hat{Y}} E[\mathcal{L}(Y, \hat{Y})] = \frac{d}{d\hat{Y}} \int_y \mathcal{L}(Y, \hat{Y}) p(y) dy$$

- Assuming interchangeability of differentiation and integration:

$$\int_y \frac{d\mathcal{L}(Y, \hat{Y})}{d\hat{Y}} p(y) dy = E[\mathcal{L}'(Y, \hat{Y})]$$

- First order condition (FOC):

$$E[\mathcal{L}'(Y, \hat{Y})] = 0$$

What is the optimal prediction under squared error loss?

Proposition

For the squared error loss, the optimal unconditional prediction is the mean of Y , $\hat{Y} = E(Y)$.

- ① Loss function $\mathcal{L}(Y, \hat{Y}) = \mathcal{L}(Y - \hat{Y}) = (Y - \hat{Y})^2$
- ② Minimize expected loss:

- First derivative:

$$\mathcal{L}'(Y, \hat{Y}) = -2(Y - \hat{Y})$$

- First order condition (FOC):

$$E[\mathcal{L}'(Y, \hat{Y})] = E[-2(Y - \hat{Y})] = 0$$

- Solution:

$$\hat{Y}^* = E[Y]$$

What is the optimal prediction under squared error loss?

This clearly extends to conditional distributions and expectations.

So for prediction we typically focus on **conditional means**, i.e. the regression function and take as predictor function

$$f(\boldsymbol{x}) = \mathrm{E}(Y | \boldsymbol{X} = \boldsymbol{x}).$$

Of course, for a new data point \boldsymbol{X} , we take $\hat{Y} = f(\boldsymbol{X})$.

Other loss functions

Mean Absolute Deviation is sometimes an alternative to MSE.¹

Proposition

For the absolute deviation loss, the optimal predictor is the median (which we assume unique).

Analogously, asymmetric linear losses leads to conditional quantiles.

Let

$$\mathcal{L}(u) = \begin{cases} -au & u < 0 \\ bu & u \geq 0 \end{cases} \quad \text{for } a, b > 0.$$

Then,

$$f(\mathbf{x}) = F_{Y|\mathbf{X}=\mathbf{x}}^{-1} \left(\frac{b}{a+b} \right)$$

Of course, $\hat{Y} = f(\mathbf{X})$ for new data \mathbf{X} .

¹Usually because of robustness to outliers; but beware of the computational burden.

Decision theory

Different losses imply different predictions for the same prediction problem!

Let $Y \sim N(0, 1)$ and I the indicator function:

$$\mathcal{L}(u, a, b, p) = (a * I_{(u < 0)} + b * I_{(u \geq 0)}) * |u|^p$$

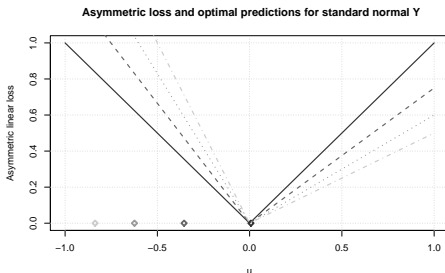


Figure: Linear loss function $p = 1$

Decision theory

Different losses imply different predictions for the same prediction problem!

Let $Y \sim N(0, 1)$ and I the indicator function:

$$\mathcal{L}(u, a, b, p) = (a * I_{(u < 0)} + b * I_{(u \geq 0)}) * |u|^p$$

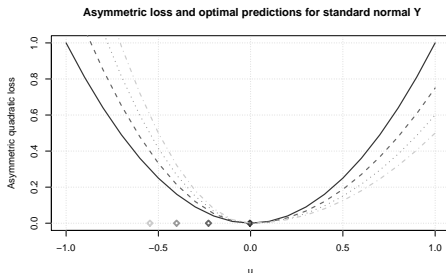


Figure: Quadratic loss function $p = 2$

Measures of location and predictions

In general, we get a specific optimal prediction for each loss \mathcal{L} .

Definition (\mathcal{L} -Measure of location)

The functional $\mu_Y^{\mathcal{L}} = \arg \min_m \mathbb{E}(\mathcal{L}(Y - m))$ is the \mathcal{L} -measure of location of Y .

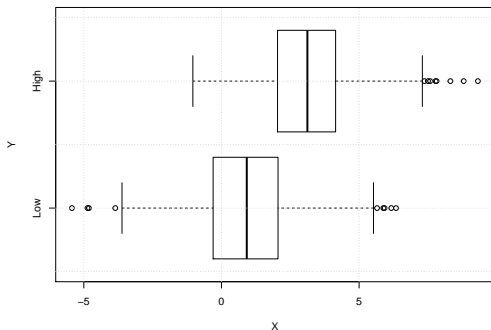
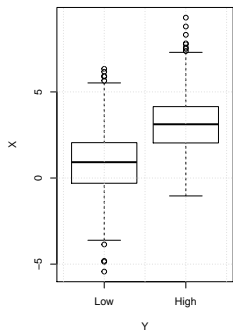
The optimal prediction \hat{Y} given \mathbf{X} under \mathcal{L} is the conditional \mathcal{L} -measure of location of Y , $\mu_{Y|\mathbf{X}}^{\mathcal{L}}$, and the predictor function is

$$f(\mathbf{x}) = \mu_{Y|\mathbf{X}=\mathbf{x}}^{\mathcal{L}}.$$

Changing the loss function changes the prediction problem!

Classification

Let's look at the binary case, with 2 possible outcomes, say 0 and 1.



Here the distribution of Y is Bernoulli, so what matters is the conditional probability that $Y = 1$ given X .

Classification errors

The so-called **confusion matrix** summarizes the error possibilities

	$\hat{Y} = 0$	$\hat{Y} = 1$
$Y = 0$	$\ell_{0,0}$	$\ell_{0,1}$
$Y = 1$	$\ell_{1,0}$	$\ell_{1,1}$

Like for prediction, you could have any cost attached to errors.

The most common loss function for binary classification is the 0/1 loss:
(We'll start with this one.)

	$\hat{Y} = 0$	$\hat{Y} = 1$
$Y = 0$	0	1
$Y = 1$	1	0

Outline

- 1 Prediction errors and quantification
- 2 Classification errors**
- 3 Learning for prediction
- 4 Up next

Optimal classifier

Proposition

The optimal binary classifier under 0/1 loss sets $\hat{Y} = 1$ if $P(Y = 1) > P(Y = 0)$ and $\hat{Y} = 0$ if $P(Y = 0) > P(Y = 1)$.

I.e. the optimal classifier

- picks the class with the larger (conditional) probability of occurrence.
- This extends to several classes, as long as one uses the 0/1 loss.

For binary classification, we may equivalently state

“Set $\hat{Y} = 1$ if conditional probability of $Y = 1$ exceeds $1/2$ ”.

Other loss functions simply imply another threshold for the conditional class probability.

Errors and the confusion matrix

Especially in machine learning, it is customary to work with error rates directly rather than losses.

We need some terminology to define them.

	$\hat{Y} = \text{No}$	$\hat{Y} = \text{Yes}$	
$Y = \text{No}$	True Negatives	False Positives	Actual No
$Y = \text{Yes}$	False Negatives	True Positives	Actual Yes
	Predicted No	Predicted Yes	Total

(These error rates are typically given in a yes/no classification situation, and we follow this convention.)

Error rates

Accuracy: $(\text{True Positives} + \text{True Negatives}) / \text{Total}$

(How often is the classifier correct?)

Misclassification/Error rate: $(\text{False Positives} + \text{False Negatives}) / \text{Total}$

(How often is the classifier wrong?)

True Positive Rate, Sensitivity or Recall: $\text{True Positives} / \text{Actual Yes}$

(How often is the classifier correct when truth is yes?)

True Negative Rate or Specificity: $\text{True Negatives} / \text{Actual No}$

(How often is the classifier correct when truth is no?)

False Positive Rate: $\text{False Positives} / \text{Actual No}$

(How often is the classifier wrong when truth is no?)

Precision: $\text{True Positives} / \text{Predicted Yes}$

(How often is the classifier correct when prediction is yes?)

Prevalence: $\text{Actual Yes} / \text{Total}$

(How often does the yes occur in the population or the sample?)

The Receiver Operating Curve

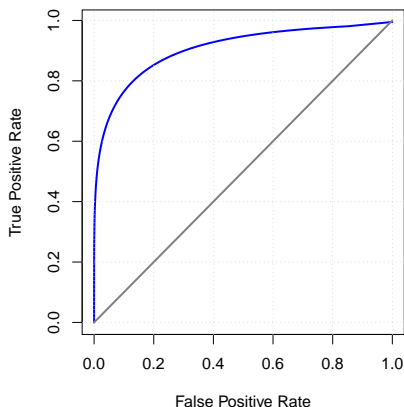
Accuracy is important, but can be misleading if prevalence is either very high or very low.

Sensitivity and specificity should be high, but there's a trade-off:

- This trade-off is controlled by the threshold (say τ) for classification based on conditional class probability.
 - lower τ means more $\hat{Y} = 1$ (higher sensitivity)
... but more of them will be false positives (lower specificity);
 - higher τ means less $\hat{Y} = 1$ & less false positives (higher specificity)
... but at the same time lower sensitivity.
- The “optimal” τ depends on the loss function ($\tau = 1/2$ for 0/1 loss).
- Plotting sensitivity (TPR) against specificity (TNR) for all $\tau \in (0, 1)$ gives the so-called Receiver Operating Curve.
- (In fact TPR is plotted for historical reasons against $FPR=1-TNR$.)

ROC & AUC

The ROC is informative about the classification performance of a model (population or fitted) for the conditional class probability.



A pure chance classifier (i.e. \hat{Y} independent of Y) has $\text{TPR} = \text{FPR}$.

So the further the ROC is away from the gray diagonal, the better!

The “Area Under the Curve” quantifies this:

- $0 < \text{AUC} < 1$
- $\text{AUC} = 1/2$ for pure chance.

Outline

- 1 Prediction errors and quantification
- 2 Classification errors
- 3 Learning for prediction**
- 4 Up next

Recall classical statistics

Say you have a parametric regression model,

$$Y = f(\mathbf{X}, \boldsymbol{\theta}) + \varepsilon$$

where $E(\varepsilon|\mathbf{X}) = 0$.

Estimate using squared error loss,

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}^*} \frac{1}{n} \sum_{i=1}^n (Y_i - f(\mathbf{X}_i, \boldsymbol{\theta}^*))^2.$$

Proposition

Regularity conditions assumed, $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}$.

Since $E(Y|\mathbf{X}) = f(\mathbf{X}, \boldsymbol{\theta})$, plug in $\hat{\boldsymbol{\theta}}$ to obtain

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}, \hat{\boldsymbol{\theta}}) \quad \text{and} \quad \hat{Y} = f(\mathbf{X}, \hat{\boldsymbol{\theta}}).$$

Consistency of $\hat{\boldsymbol{\theta}}$ implies that $\hat{f}(\mathbf{x}) \xrightarrow{p} f(\mathbf{x})$ and $\hat{Y} \xrightarrow{p} E(Y|\mathbf{X})$.

Prediction errors

The prediction errors at a (new) data point (Y, \mathbf{X}) are

$$Y - \hat{Y} = \varepsilon - \left(f(\mathbf{X}, \hat{\boldsymbol{\theta}}) - f(\mathbf{X}, \boldsymbol{\theta}) \right).$$

Since $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}$, we may linearize f in $\boldsymbol{\theta}$ (at the given \mathbf{X}) and obtain

$$\hat{\varepsilon} = \varepsilon - \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \right)' \nabla f + O_p\left(\frac{1}{n}\right)$$

where ∇f is the gradient of f at the population $\boldsymbol{\theta}$ and the given \mathbf{X} .

We thus have **estimation error** and **irreducible error**.

The estimation error can be made arbitrarily small if there are enough data.

Misspecification

Say you want to fit the same parametric $f(\mathbf{x}, \boldsymbol{\theta})$, but the actual regression function is different, say $E(Y|\mathbf{X}) = m(\mathbf{X})$.

Proposition

Let $\boldsymbol{\theta}_p = \arg \min_{\boldsymbol{\theta}^*} E \left((f(\mathbf{X}, \boldsymbol{\theta}^*) - m(\mathbf{X}))^2 \right)$ (we call $\boldsymbol{\theta}_p$ pseudo-true value). Regularity conditions assumed, $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_p$.

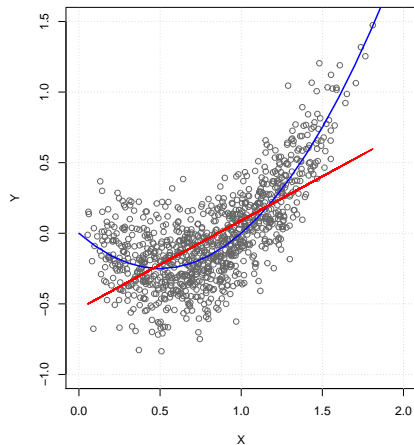
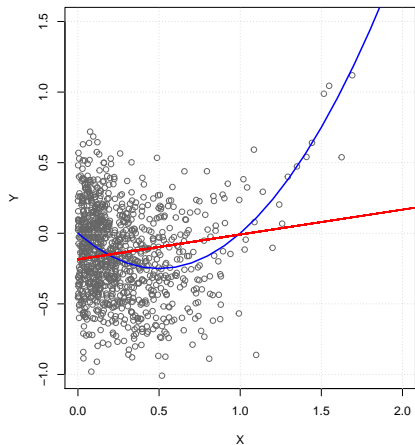
Regarding prediction errors,

$$Y - \hat{Y} \approx \varepsilon + (m(\mathbf{X}) - f(\mathbf{X}, \boldsymbol{\theta}_p)) - (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_p)' \nabla f$$

where ∇f is the gradient of f at \mathbf{X} for the pseudo-true $\boldsymbol{\theta}_p$.

We thus have **model bias** as an additional source of error.

Example: Fitting linear regression instead of quadratic



Note that the bias depends on the marginal distribution of X !

Sideline remark: Empirical risk minimization

It is important to realize that you get what you fit!

Say you use a different loss function,

$$\hat{\theta} = \arg \min_{\theta^*} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(Y_i - f(\mathbf{X}_i, \theta^*)).$$

Proposition

Regularity conditions assumed, $f(\mathbf{x}, \hat{\theta}) \xrightarrow{p} \mu_{Y|\mathbf{X}=\mathbf{x}}^{\mathcal{L}}$.

E.g. fitting under MAD gives the (fitted) conditional median as predictor.

So you should typically fit predictive models using the same loss function which is used to evaluate the errors.

Flexibility I

Some parametric families of functions are **universal approximators**.

I.e., as $\dim \boldsymbol{\theta} \rightarrow \infty$,

- Polynomials, Fourier series, splines, and also
- regression trees and neural networks,

can approximate smooth functions arbitrarily well.

So choosing “large” $\dim \boldsymbol{\theta}$

- gives enough flexibility of the model such that it
- accommodates arbitrary unknown regression functions.
- The analogous findings holds for classifiers.

Free lunch anybody?

Flexibility II

The problem is that

$$\text{Var}(\hat{\varepsilon}) \approx \text{Var}(\varepsilon) + \nabla f' \text{Cov}(\hat{\theta}) \nabla f$$

where

$$\nabla f' \text{Cov}(\hat{\theta}) \nabla f \leq \|\nabla f\|^2 \|\text{Cov}(\hat{\theta})\|.$$

In the worst-case scenario, $\|\text{Cov}(\hat{\theta})\| = O(\dim \theta)!$

You buy bias reduction by increasing the variance!

The bias-variance trade-off² is really everywhere...

²Further readings: ESL Chapter 7.2 Bias, Variance and Model Complexity

Bias-Variance trade-off

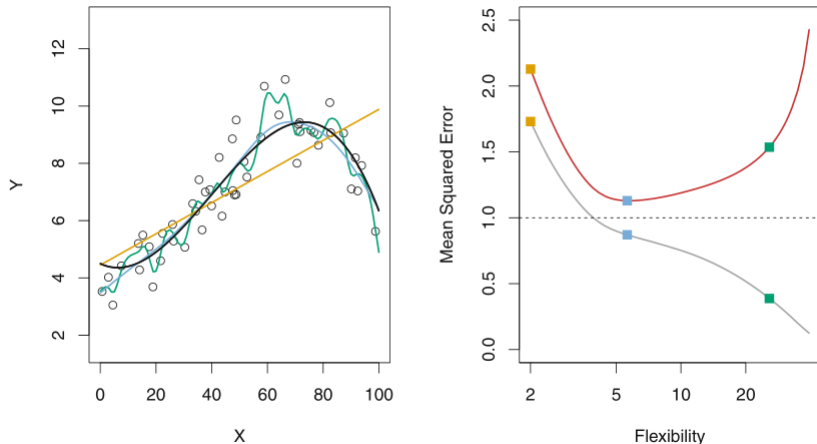


Figure: Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel. FIGURE 2.9. ISLR

Towards model selection

One minimizes prediction MSE by balancing bias and estimation error.

- This is a critical step, since both influence prediction MSE,
- and can be seen as a model selection issue.

In statistical learning we have new names:

- Bias dominates: **underfitting** (oversmoothing)
- Variance dominates: **overfitting** (undersmoothing)

Choosing the right amount of smoothing is sometimes method-specific so we deal with it on suitable occasions.

More information is a good thing, right?

Intuitively,

The more features you have the better!

E.g. chances increase that a powerful predictor is among them.

- You need however enough data to be able to incorporate the extra features in your model.
- Clearly, additional features (and correspondingly parameters) for the same amount of data increase estimation variance.

The curse of dimensionality

The **curse of dimensionality** is a generic term used to describe situations where it is very costly to exploit additional features.

It originated in dynamic optimization, where the cost of adding one dimension to a grid search increases exponentially with the number of dimensions.

Its exact form depends on the model, but the source is the combination of high dimensionality with high flexibility.

The curse is not only for prediction!

Classification models also suffer from the curse of dimensionality.

And unsupervised learning has its dimensionality issues as well.

- Imagine a clustering task with n data points and $p = n$ features.
- Say for simplicity that the n features are binary 0/1,
- ... and $X_{i,j} = 1$ for $i = j$ only.
- Unless you pay attention, most algorithms will immediately find n different clusters

For all, **dimensionality reduction** may help deal with the curse.
(And so does **model selection**.)

Outline

- 1 Prediction errors and quantification
- 2 Classification errors
- 3 Learning for prediction
- 4 Up next**

Coming up

Linear prediction