ft → TRAINING SET $(x_1,y_1),...,(x_n,y_n)$ (independent or)
TEST SET: Same from
model assess → VALIDATION

VS. MODEL ASSESMENT

## LINEAR REGRESSION (M3)

classical normal regression model
$$Y = X\beta + \varepsilon$$
$$\varepsilon \sim N_n(0, \sigma^2 I)$$

$X$ here, rank $p+1$ (full rank) $n > p+1$

$\hat{\beta} = (X^TX)^{-1}X^TY$
$\hat{\beta} \sim N_{p+1}(\beta, \sigma^2(X^TX)^{-1})$

$\hat{\sigma}^2 = \frac{RSS}{n-p-1}$

$c_{jj}$ is the diagonal element of $(X^TX)^{-1}$

Inference for one $\beta_j$: based on
$$T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}}\hat{\sigma}} \sim t_{n-p-1}$$

Hypothesis test:
$H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$
Reject $H_0$ when $T_0 = \frac{\hat{\beta}_j - 0}{\sqrt{c_{jj}}\hat{\sigma}}$

P-value: $P(C > |t_{obs}|)$, $t_{obs}$ observed F.

residuals: $e_i = y_i - \hat{y}_i$ and
Scaled versions thereof → plot

### [M2+3] BIAS-VARIANCE trade-off in regression setting: $Y = f(x) + \varepsilon$

Expected test mean squared error at $x_0$
$$E[(Y - \hat{f}(x_0))^2] = \cdots = Var(\varepsilon) + Var(\hat{f}(x_0)) + [Bias(\hat{f}(x_0))]^2$$
irreducible error

PLOT OF THIS:

VERY POPULAR PLOT OF TRAINING & TEST ERROR

minimize 0/1 loss

## CLASSIFICATION (M4)

$(Y_i, x_i)$ with $Y_i \in \{1,...,K\}$ or $Y_i = \{0,1\}$
$p_k(x_0) = P(Y=k | X=x_0)$ posterior class probability
$\pi_k = P(Y=k)$ prior class probability
$f_k(x) = P(X=x | Y=k)$ pdf of $X$ in class $k$ class conditional distribution

### Sampling paradigm: LDA & QDA

LDA: $f_k(x)$ is $N(\mu_k, \Sigma)$, $\hat{\mu}_k = \frac{1}{n_k}\sum_{i:y_i=k} x_i$
QDA: $f_k(x)$ is $N(\mu_k, \Sigma_k)$, $\hat{\Sigma}_k = \frac{1}{n_k-1}\sum_{i:y_i=k}(x_i-\hat{\mu}_k)(x_i-\hat{\mu}_k)^T$
$\hat{\pi}_k = \frac{n_k}{n}$

$\pi_k f_k(x) = \pi_j f_j(x)$ → linear class boundaries for LDA, quadratic

### Diagnostic paradigm: KNN & logistic regression (trees & SVM too)

KNN: $\hat{P}(Y=j | X=x_0) = \frac{1}{K}\sum_{i\in N_0} I(y_i=j)$
$N_0$ neighborhood of $x_0$

Logistic regression $(k=2)$: $Y_i = \{0,1\}$ with $p(1-p)$
where $p_i = \frac{\exp\{x_i^T\beta\}}{1+\exp\{x_i^T\beta\}}$

logistic function also sigmoid

Parameter estimates found by maximum likelihood (TMA4300 em) not on closed form.

odds: $\frac{p}{1-p}$ so $p_i = 1$ gives odds 1, and is relevant

Class boundary is linear in $x$-space, because
$\log\left(\frac{P(Y=1)}{P(Y=0)}\right) = x^T\beta$ and set a cut-off on $p$

### CONNECTION SVM & logistic regression

SVM: $\min_{\beta} \left\{\frac{1}{n}\sum_i \max(0, 1-y_i f(x_i)) + \lambda\sum \beta_j^2\right\}$
hinge loss, ridge penalty

### SUPPORT VECTOR MACHINES (M9)

A method both for regression and classification, but we only consider classification, and two classes

Aim: find hyperplane that separate (perfectly) the two classes
$\beta_0 + x^T\beta = 0$ with normalized $\beta$, $\sum_{j=1}^p \beta_j^2 = 1$

$\beta_0 + x^T\beta > 0$ one side of hyperplane
$< 0$ the other side

$y_i \cdot (\beta_0 + x_i^T\beta) > 0$ if correct classified $x_i$

Maximal margin classifier:
$\max_{\beta_0,...,\beta_p} M$ subject to $\sum_{j=1}^p \beta_j^2 = 1$
$y_i(\beta_0 + x_i^T\beta) \geq M$ $\forall i=1,...,n$

Support vector classifier: non-separable case — $c_0, \varepsilon_i$: slack variables
$\max_{\beta_0,...,\beta_p,\varepsilon_1,...,\varepsilon_n} M$ subject to $\sum_{j=1}^p \beta_j^2 = 1$, $y_i(\beta_0 + x_i^T\beta) \geq M(1-\varepsilon_i)$ $\forall i$, $\varepsilon_i \geq 0$, $\sum_{i=1}^n \varepsilon_i \leq C$

C is tuning parameter chosen by CV

Classification rule: $f(x^*) = \beta_0 + x^T\beta$ and if $f(x^*) < 0$ set $y^* = -1$, $> 0$ set $y^* = 1$

Support vector machine: $f(x) = \beta_0 + \sum_{i\in S}\alpha_i K(x,x_i)$
$x_i^T\beta$ replaced by more general nonlinear functions that are written as sums of kernels

popular kernel: radial kernel $K(x_i,x_{i'}) = \exp\left(-\gamma\sum_{j=1}^p (x_{ij}-x_{i'j})^2\right)$

tuning: $(C, \gamma)$ by CV

### UNSUPERVISED LEARNING (M10)

Look for underlying structure or groupings in data — no $Y$ only $X_{n\times p}$

Principal component analysis (see M6)
PC loadings: interpret effect of each covariate on each component
PC score: express effect, diff between A and B etc

### MODEL SELECTION (M6)

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2, \quad \bar{y} = \frac{1}{n}\sum y_i$$ total sum of squares

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$ residual sum of squares

$R^2_{adj} = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)}$

AIC = $-2\log(\text{likelihood}) + 2(\text{number of param})$
BIC = $-2\log(\text{likelihood}) + \log(n) \cdots$

penalizing of training error

$$R^2 = \frac{TSS-RSS}{TSS} = 1 - \frac{RSS}{TSS}$$
Coeff of determination proportion of variance explained by regression

RSS will always decrease (or stay unchanged) if new covariates are added to the regression (even if they are just noise), and $R^2$ will always increase. Therefore $R^2$ can never be a criterion to choose between eg models of different model complexity based on a training set alone.

### DIMENSIONALITY REDUCTION: Principal component regression - PCA (M6&10)

Have $X$ with $p$ large and/or many correlated covariates: want set of linear combs of the $X$s that captures large part of the variability.

Let $\Sigma$ be covariance (or correlation) matrix of $X$, and $(\alpha_i, \lambda_i)$: its eigenvalue/eigenvector pairs of $\Sigma$ where $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$. Then construct $Z_i = \alpha_i^T X$.

How many PC's?
— proportion of total variance explained by first $k$ PC's $\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^p \lambda_i}$

### LOSS function

regression: Quadratic: $\sum_{i=1}^n (y_i - f(x_i))^2$

classification 0/1: $y_i = \hat{y}_i$ 0 loss, $y_i \neq \hat{y}_i$ 1 loss

Bayes classifier: classify to the class with the highest probability $P(Y=k | X=x)$ — and in simulation experiments this is known.

Bayes class boundary: boundary of Bayes classifier.
Bayes error: error rate of Bayes classifier.

### RESAMPLING METHODS (M5)

DATA rich situation (sometimes, e.g. when we generate data ourselves)

| TRAIN | VALIDATE | TEST |
|-------|----------|------|
| fit model | model selection (complexity) | model assessment |

DATA not so rich: Common situation
→ Validation set approach
→ Leave-one-out cross-validation: LOOCV
→ k-fold cross-validation (CV)

Assess with MSE: $(y_i - \hat{y}_i)^2$ for each obs → regression; or 0/1: $I(y_i \neq \hat{y}_i)$ → classification. Sum over folds and average

How to PERFORM CV: ALL PARTS of model fitting must be in the CV loop unless → selection bias. Eg cancer predicted perfectly from J. genes example.

### MORE RESAMPLING (M5)

BOOTSTRAPPING: data set of $n$ observations. Estimate $f$ key $t$.
Draw random sample from $f$ → draws with replacement from data → one bootstrap sample

The probability that $i$ is in bootstrap sample is
$1 - (1-\frac{1}{n})^n \to 0.632$ ($n\to\infty$)

We may estimate $SD(\hat{\beta})$ from regression — easy bootstrapping. CI's in TMA4300 Comp stat.
$\hat{SD}(\hat{\beta}) = \sqrt{\frac{1}{B-1}\sum_{b=1}^B(\hat{\beta}_b - \frac{1}{B}\sum_{r=1}^B \hat{\beta}_r)^2}$
$B$: total no of bootstrap samples

### Assessment of classification

Confusion matrix:

|  | | predicted class | |
|--|--|--|--|
| | | 0 | 1 |
| true class | 0 | TN | FP | N |
| | 1 | FN | TP | P |

may report misclassification rate: $1 - \frac{TN+TP}{n}$
accuracy: rate of correct classifications $\frac{TN+TP}{n}$

Sensitivity: $TP/(PN+TP)$
Specificity: $TN/(TN+FP)$

ROC-curve, each is based on one cut-off on $p$.
AUC: area under ROC-curve
AUC = 1: random guessing

### MODEL REGULARIZATION (M6)

Motivation: LS regression gives unbiased estimate but variance might be large, e.g. when covariates are correlated → $(X^TX)^{-1}$. If we instead look at biased $\hat{\beta}$ we might get smaller $Var(X^T\hat{\beta})$.

### Ridge (L2) regression:

$\min_{\beta} RSS + \lambda\sum_{j=1}^p \beta_j^2$
$\lambda$: tuning parameter decided by CV value giving min MSE
not scale invariant - useful to scale covariates before using ridge.

### Lasso (L1) regression:

$\min_{\beta} RSS + \lambda\sum_{j=1}^p |\beta_j|$ (automatically performs model selection)

Dual problem:
$\min$ RSS subject to $\sum_{j=1}^p |\beta_j| \leq s$ lasso
$\sum_{j=1}^p \beta_j^2 \leq s$ ridge

### Partial least squares (PLS): Competitor to PCR, but now $Z$ is found also taking $Y$ into account. More details in TMA4267.

### MOVING BEYOND LINEARITY (M7)

$Y = f(x) + \varepsilon$ and MLR $f(x) = X\beta$ → $\hat{\beta} = (X^TX)^{-1}X^TY$
$X$ can be replaced by a set of basis functions $\tilde{X} = \begin{bmatrix} b_1(x) & b_2(x) & ... \end{bmatrix}$
So we still have linearity in $\beta$'s — but now $Y = \tilde{X}\beta$, so $\hat{\beta} = (\tilde{X}^T\tilde{X})^{-1}\tilde{X}^TY$

1) POLYNOMIAL regression: $b_k(x) = x^k$

2) STEP-FUNCTIONS: $b_k(x) = I(C_k \leq x < C_{k+1})$
similar to dummy variable coding, but now $x$ might be continuous.

3) REGRESSION SPLINES: combine polynomials & steps at knots
order-1 spline: $f(x) = x^{j-1}$, $j=1,...,M-1$
$b_{k+M-1}(x) = (x - C_k)_+^{M-1}$; $k = 1,..., K$
where knots $C_1,...,C_k$ are chosen (or estimated)
The resulting curve is smooth and deriv. up to order $M-2$ continuous (also in knots). Cubic spline popular.

NATURAL CUBIC SPLINE: cubic spline that is linear at ends

$\lambda$ can be chosen by CV, or some version of AIC. Then effective number of parameters is calculated from the trace of the smoother matrix. (beyond our scope)

4) SMOOTHING SPLINE - penalization problem (similar to ridge & lasso) $g(x)$ smoothy spline
$$\min \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$
$\lambda = 0$ interpolates ... penalize variation
It turns out (magic) that a smoothing spline is a natural cubic spline with knots at the unique values of $x$.

5) LOCAL regression: smoothed version of KNN regression
$\min \sum_{i=1}^n K_{i0}(y_i - \beta_0 - \beta_1 x_i - f \cdot x_i)^2$
close observation higher weight than far apart. different versions

6) Additive model: put it all together
$\tilde{X} = [1 \; f_1(x_1) \; f_2(x_2) \; ...]$
EXTENSIONS: this additive model can be used on $p_i = \frac{\exp(x_i^T\beta)}{1+\exp(x_i^T\beta)}$ to generalize to logistic regression make then non-linear in $X$-space.

### TREE-BASED METHODS (M8)

Both for regression & classification — and we have nonlinearities & interactions between covariates!

1) Want to minimize an error criterion on the whole tree construction — but since computationally infeasible a greedy approach "recursive binary splitting" is used.

$R_1(j,s) = \{x | x_j < s\}$, $R_2(j,s) = \{x | x_j \geq s\}$: find predictor $(j)$ and splitting point $(s)$ that minimize "RSS"
regression: $\sum_{i:x_i\in R_1}(y_i - \hat{y}_{R_1})^2 + \sum_{i:x_i\in R_2}(y_i - \hat{y}_{R_2})^2$
$\hat{y}_{R_1}$: mean of $Y$ for training obs in $R_1$
$\hat{p}_k = \frac{1}{N_k}\sum_{i\in R_k} I(y_i=k)$
classification: similar, but using impurity measure:
• Gini index: in region $R_j$ (node) with $N_j$ obs.
$G = \sum_{k=1}^K \hat{p}_{jk}(1-\hat{p}_{jk})$
• Cross entropy: $D = -\sum_{k=1}^K \hat{p}_{jk}\log(\hat{p}_{jk})$
• Deviance: $-2\sum_k \hat{p}_{jk}\log(\hat{p}_{jk})$

2) Continue splitting until no eg: reduction in RSS not big, less than $\delta$ obs in a terminal node.

3) This full tree $T$ used for prediction by dropping a new obs down the tree and predict as
a) regression: $\hat{y} = \hat{y}_{R_j}$ = $\frac{1}{N_j}\sum_{i:x_i\in R_j} y_i$
b) classification: $\hat{y}_j = $ majority vote or $\hat{p}_k = \frac{1}{N_j}\sum_{i:x_i\in R_j} I(y_i=k)$ and classify to the class with the highest probability

4) From a full tree to a pruned tree: reduce the size of the tree by cutting branches.
Cost complexity pruning: for a complexity param $\alpha$ minimize $\sum_{T} Q(T) + \alpha|T|$
$|T|$: number of terminal nodes; $\alpha$ found by CV.

5) Evaluate model fit on test set. Same criteria as regression (RSS) or classification (ROC, AUC, misclassification rate).

### A) BAGGING

Bootstrap aggregation
Motivation: $Var(\bar{x}) = \frac{\sigma^2}{n}$ → we construct many trees and average them to define $\hat{y}$. The $B$ trees we grow (full) from $B$ bootstrap samples $\hat{f}_{avg}(x) = \frac{1}{B}\sum_{b=1}^B \hat{f}_b(x)$

### B) RANDOM FOREST — but not maximal effect of bagging

Because trees are not independent → we make them more uncorrelated by not allowing all $p$ covariates to be chosen in each binary split.
Rule of thumb: $m = \sqrt{p}$ classification and $m = \frac{p}{3}$ regression.

OOB - out-of-bag: for each bootstrap sample there are on average $\frac{2}{3}$ of the observations in the sample (bag). The last $\frac{1}{3}$ are used for prediction. We will then have ~$\frac{B}{3}$ predictions for observation $\ell$ — which we average.

Variable importance plots: or Gini
— give the total amount of RSS decrease over splits of a predictor — averaged over all $B$ trees.

### C) BOOSTING: many repeats, but we only considered boosting of trees

Fit one tree — make residuals — fit a tree to the residuals → update
$\hat{f}(x) = \sum_{k=1}^B \lambda \hat{f}_b(x)$
$r \leftarrow r - \lambda \hat{f}_b(x)$

Tuning parameters: $\lambda$: shrinkage parameter; $d$: number of tree splits; $B$: number of trees

### NEURAL NETWORKS (M11)

a) Possible to represent MLR & logistic regression as graph (one input and one output layer)
$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p + \varepsilon$
$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + ... + \beta_p x_p$, $p = \frac{\exp(\alpha_{c-1})}{1+\exp(\alpha_c)} = \frac{\exp(\alpha_c)}{\sum \exp(\alpha_k)}$
for $K$ classes: $K$ output (dummy) and softmax activation
used for two-class classification problems

b) Instead of MLR where we add nonlinear function of each covariate, we instead look at nonlinear function of sums of covariates pr. layer — and many layers. Popular non-linear activation function is relu(a) = max(0,a). For each layer we specify number of nodes and choice of activation function → feedforward network.

technical issue: data represented as tensors (generalization of vector, matrix)

c) How to find/estimate the unknown weight (= parameters)?
— minimize $\sum$ loss: quadratic loss for regression "mse"
$K=2$ classes: "binary crossentropy" = as in tree $-\frac{1}{n}\sum_i (y_i\log(p_i) + (1-y_i)\log(1-p_i))$
$K$ classes: "categorical crossentropy" $-\frac{1}{n}\sum_i\sum_k y_{ik}\log(p_{ik})$
— optimization: backpropagation, can use chain rule of differentiation
gradient based: popular → mini-batch stochastic gradient descent with stochastic version with "batches" of observations
epoch: one run through all observations

d) Evaluation of performance: "metrics"
mse = mean abs. error, accuracy: avg correct classif.

e) avoiding overfitting:
— smaller network or more data
— regularization (L1, L2 or dropout)
— early stopping

f) "unsolved": how do we assess the uncertainty in the network fitted?

g) many cool extensions
we have briefly looked at recurrent nets & convolution nets

h) data types popular: images, text that M11 do not have solved.

### Hierarchical clustering

— work in a sequential way by connecting observations that are similar

### K-means clustering (non-overlapping clusters)

* Number of clusters given (or selected separately)
* Define cluster centroids = densify observations to clusters (closest)
* Recalculate centroids
Euclidean distance popular (but also Mahalanobis)
Generalizations to medoids etc.

Similarity measures: Euclidean (dissimilarity), Correlation
weighted
Linkage: how to calculate dissimilarity between groups of observations?
— single: minimum
— average: average
— complete: maximum
* Presented in a dendrogram — choose where to cut the dendrogram to get a number of clusters.