

DD2421 Machine Learning

K-nearest neighbor

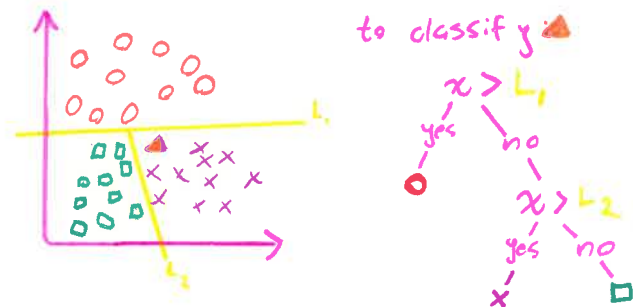
Presented with an unlabeled sample its predicted label will be the same as the majority of the k nearest neighbor.



If k increases the decision boundary becomes smoother, it also generalizes better.
Memory & computational requirements scale linearly with training size.
Lower k higher overfitting

Decision Tree

Create logical questions that partitions the search space into classes.



train by finding best questions that maximize information gain, and splitting data sets accordingly.

Information gain is measured in entropy.

$$\text{Entropy} = \sum_i -p_i \log_2 p_i \quad \forall p$$

1 average 2 binary

Entropy is the number of questions needed to find an answer.
aka number of bits of information

Alternative: Gini Impurity

$$\sum p_i(1-p_i) = 1 - \sum p_i^2$$

Expected error rate at node N if label from class distribution at N similar to entropy but peaks stronger at equal probabilities.

Idea:

Asks question that maximizes expected reduction of entropy

$$\text{Gain} = \underbrace{\text{Entropy}(S)}_{\text{Before}} - \sum_S \underbrace{\frac{|S_i|}{|S|}}_{\text{weighted subset}} \underbrace{\text{Entropy}(S_i)}_{\text{entropy of subset}}$$

After

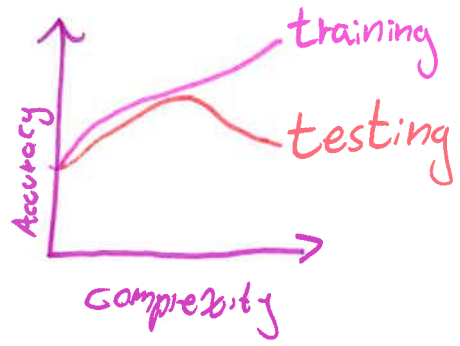
Choose attribute that says the most about the answer.

Overfitting

The model performs too specialized on training data

↳ bad representation ↳ too complex model ↳ noisy examples

Three datasets, training, validation, testing



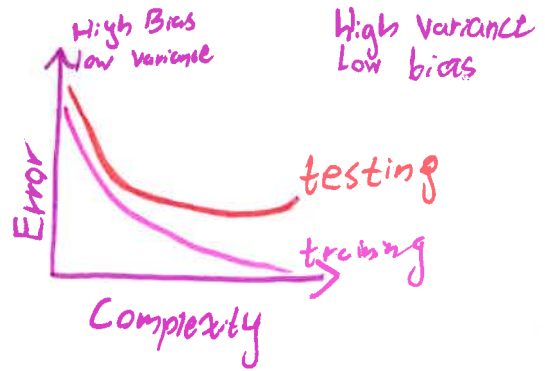
K-fold Validation:

Split data into k folds

train k times using different validation set

Bias: Difference between average prediction and correct value

Variance: Variability in model prediction between realizations of model



Regression

fit a function to a dataset

$$D = \{(\bar{x}_1, \bar{y}_1), \dots, (\bar{x}_n, \bar{y}_n)\}$$

Linear regression

$$\hat{f}(\bar{x}) \approx \sum w_i x_i = W^T \bar{x}$$

Mean squared error

$$(\hat{f}(\bar{x}) - f(\bar{x}))^2$$

$$E(\hat{f}) = \frac{1}{N} \sum (\hat{f}(\bar{x}_n) - f(\bar{x}_n))^2$$

$$= \frac{1}{N} \| \bar{X}W - Y \|^2$$

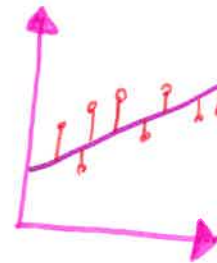
compute \hat{w} that minimizes E

Residual Sum of Squares

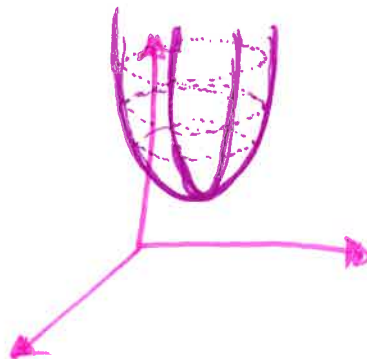
$$E_{in}(w) = \| Xw - Y \|^2$$

$$\frac{\partial}{\partial w} E_{in}(w) = 2X^T(Xw - Y)$$

When gradient is zero the error is minimized



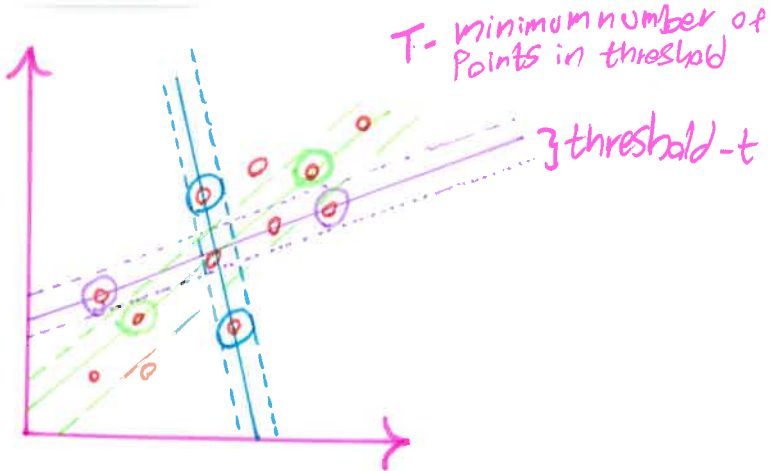
Least squares



Random Sampling Consensus

Ransac, more robust regression

- 1) sample 2 points
- 2) Draw line with a threshold
- 3) repeat from 1 until line is good enough.



2 hyperparameters:
 distance threshold - t
 subset threshold - T

Least Absolute Shrinkage & Selector Operator

Lasso Regularization aka l_1 norm. Will force some coefficients to zero
 This is added as a penalty to the loss function

$$\lambda \sum |w_i|$$

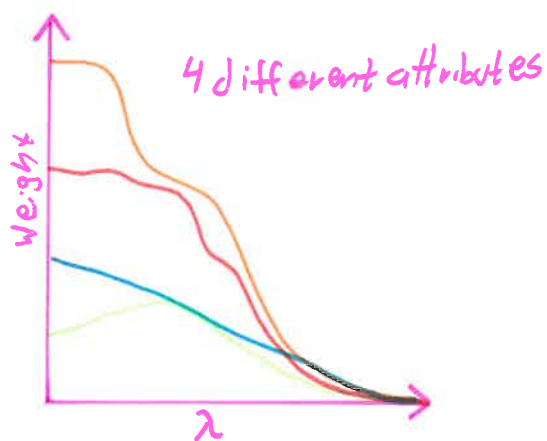
\uparrow factor \uparrow total weights

λ increases Bias while decreasing Variance

Ridge Regression

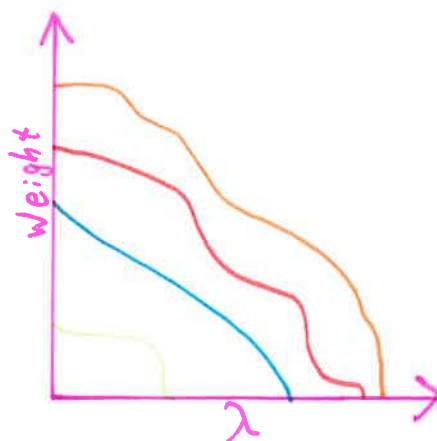
Similar to Lasso but squared

$$\lambda \sum w_i^2$$



All weights for the attributes tend to zero as $\lambda \rightarrow \infty$ but they are always above zero

Ridge Regression



All weights decrease and eventually reach zero. Some do it before others

Lasso

Probabilistic Learning

axiom $P > 0, \sum P = 1$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \rightarrow \text{if } A \text{ and } B \text{ are independent}$$

$$P(A \cap B) = P(A)P(B)$$

$$\rightarrow P(A|B) = P(A)$$

Bayes Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

\swarrow Posterior \swarrow Likelihood \downarrow evidence \swarrow Prior

Probabilistic Regression

1) find Joint distribution of X and Y

$$P(x, y)$$

2) compute posterior of Y

$$P(y|x)$$

3) compute Expectation

$$E(y|x)$$

Probabilistic Classification

Same as regression, but use Argmax

Learning As Inference

Discriminative & Generative models

\hookrightarrow Discriminates classes

$$\hookrightarrow P(y|x, D)$$

\hookrightarrow Draw boundary

\hookrightarrow tries to model

underlying distribution

$$\hookrightarrow P(x|y, D)$$

\hookrightarrow Draw distribution



Parametric vs Nonparametric

$$P_r(y|x) = P(y|x, \theta)$$

- \rightarrow Estimate $\hat{\theta}$ using D
- \rightarrow compute $P(y|x, \hat{\theta})$
- \rightarrow learning means finding $\hat{\theta}$

Estimate $P(\theta|D)$

Compute $P(y|x, D)$ from $P(y|x, \theta, D)P(\theta|D)$ by marginalizing out θ

Number of parameters grow with data

Fundamental Assumptions

\hookrightarrow observations are independent and identically distributed

Thus: $D = \{o_1, o_2, \dots, o_N\}$ $o_i = (x_i, y_i)$

$$P(D) = \prod_i P(o_i)$$

Take log

$$\log(P(D)) = \sum_i \log(P(o_i))$$

Maximum Likelihood Estimate

Find parameters that make D most likely.

$$\Theta_{ML} = \underset{\Theta}{\operatorname{argmax}} P(D|\Theta) = \underset{\Theta}{\operatorname{argmax}} \log(P(D|\Theta))$$

This can be used to approximate the x or y

$$P(x|y, D) \approx P(x|y, \Theta_{ML}) \quad P(y|x, D) \approx P(y|x, \Theta_{ML})$$

Naive Bayes Classifier

Assumes features are independent.

$$\operatorname{MAP}_{NB} = \underset{y}{\operatorname{Argmax}} P(y) \cdot \prod_D P(x_d|y)$$

Generative model
 $P(x|y)$

Logistic Regression

Treating a binary classification problem as regression problem.

Discriminative
 $P(y|x)$

$$Y|X \sim \text{Bernouli}(\lambda(x))$$

$$P(y|x) = \lambda(x)^y (1 - \lambda(x))^{(1-y)}$$

$$\lambda(x) = \text{sigmoid}(w^T x)$$

Maximum A Posteriori

Choose most likely Θ given D

$$\Theta_{MAP} = \underset{\Theta}{\operatorname{argmax}} P(\Theta|D)$$

$$= \underset{\Theta}{\operatorname{argmax}} \frac{P(\Theta) P(D|\Theta)}{P(D)}$$

$$= \underset{\Theta}{\operatorname{argmax}} P(\Theta) P(D|\Theta)$$

$$P(\Theta) \cdot \prod_N P(x_n|\Theta)$$

$$\underbrace{\log(P(\Theta))}_{\text{works as a regularizer}} + \sum_N \log(P(x_n|\Theta))$$

Bayesian Estimation

- 1) consider Θ as a random variable
- 2) Characterize Θ with $P(\Theta|D)$
- 3) Compute new posterior $P(y|x, D)$
marginalize over Θ

$$P(y|x, D) = \int_{\Theta \in \Theta} P(y|x, \Theta) P(\Theta|D) d\Theta$$

Clustering K-Means

- 1) Place k cluster centers in D
- 2) assign each $d \in D$ to closest cluster.
- 3) Move cluster centre to the mean position of all $d \in$ cluster
- 4) repeat 2-3 until mean difference for points is less than ϵ

Expectation Maximization

Mix of K-means and GMM.

- 1) for each sample introduce h_i
- 2) estimate which gaussian distributed model is most responsible for x_i
 $P(h|x)$ ← called responsibility

Place K random gaussian
for each point do $P(h|x)$
update gaussians accordingly

ML	D	Θ_{ML}	$P(y x, \Theta_{ML})$
MAP	$D, P(\Theta)$	Θ_{MAP}	$P(y x, \Theta_{MAP})$
BAYES	$D, P(\Theta)$	$P(\Theta D)$	$P(y x, D)$

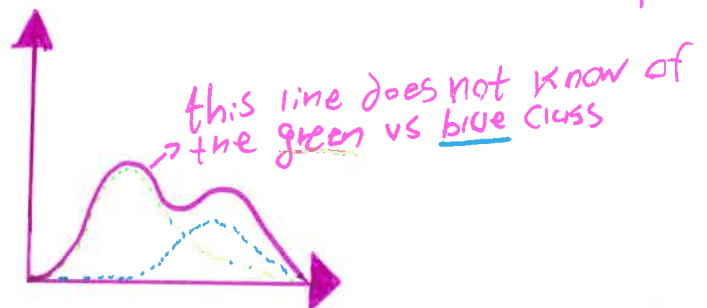
Gaussian Mixture Models

distribution is a weighted sum of gaussian distributions.

$$P(x|\Theta) = \sum_{k \in K} \tilde{\pi}_k P(x|\Theta_k)$$

$$\Theta = \{ \tilde{\pi}_1, \tilde{\pi}_2, \tilde{\pi}_3, \dots, \tilde{\pi}_K, \Theta_1, \Theta_2, \dots, \Theta_K \}$$

We cannot assume classes are independent



Linear Separation

Use a line to separate data into True and False. Train Weights and bias to correctly separate classes

In high-dimensional spaces linear separation becomes easy.

Perceptron Learning

Only change weights when input is wrong. will always find solution if it exists.

$$W_i = W_i + \eta(t - o)x_i$$

May give bad generalization since first viable solution is returned.

Risk Minimization

find Plane that has the largest margin to datapoints

$$W^T x = 0 \text{ any point in input}$$

$$t \vec{w}^T \vec{x} \geq M \rightarrow \text{margin}$$

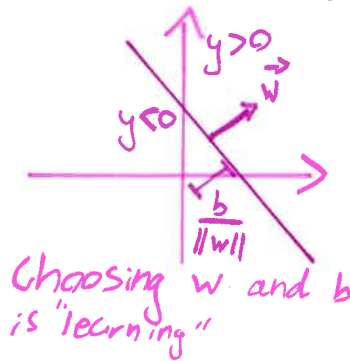
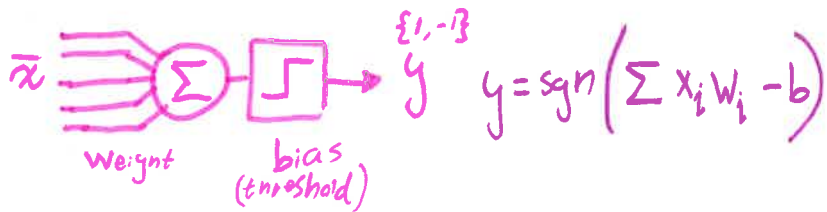
+1 for positive target
-1 for negative target

Best margin is when

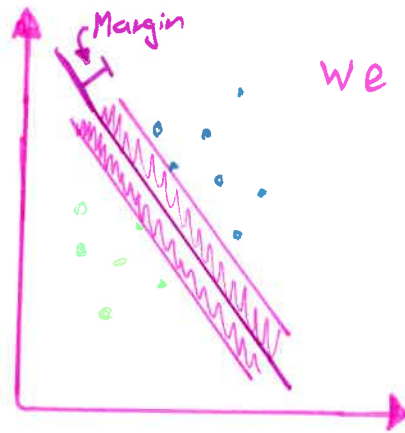
$$\frac{1}{\|w\|} \text{ is maximized.}$$

$$\Rightarrow \|w\| \text{ is minimized}$$

Artificial Neuron



Always a straight line. Does not work on non-linearly separable data.



We want to maximize M

SUPPORT Vector Machine

Project data into higher dimension
Projections use kernel to prevent there is a projection while saving resources.

Kernel:

$$\phi(\vec{x})^T \cdot \phi(\vec{y}) = k(\vec{x}, \vec{y})$$

New equation

$$t \vec{w}^T \phi(x) \geq M$$

Slacks:
introduce
 $C \sum \xi_i$

$$t \vec{w}^T \phi(x) \geq M + \xi_i$$

Lagrange multiplier

$$L = \frac{1}{2} W^T W - \sum \alpha_i [t_i w \phi - 1]$$

minimize L w.r.t w

$$\Rightarrow W = \sum_i \alpha_i t_i \phi(x_i)$$

$$\Rightarrow \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j t_i t_j \underbrace{\phi(x_i) \phi(x_j)}_{\text{kernel}}$$

some will tend to zero?

Principal Component Analysis

Reduce the number of variables by only dealing with components with large variance.

Premise: you are given many factors which describe a dataset but some are correlated meaning they can be combined.

Projecting multiple dimensions onto a single vector \vec{u} will reduce dimensions.

$$x' = (\vec{x}^T \cdot \vec{u}) \vec{u}$$

information preserved

$$(\vec{x}^T \cdot \vec{u})^2 \rightarrow \text{maximize}$$

eigenvalue

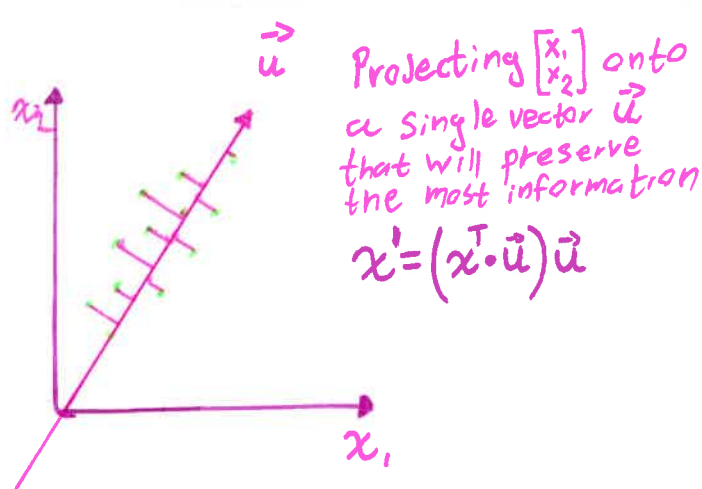
$$C\vec{u} = \lambda\vec{u}$$

eigenvector

since C is symmetric, all \vec{u} will be orthogonal

The eigenvalue determines the weight of the component.

Covariance matrix



Projecting $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ onto a single vector \vec{u} that will preserve the most information
 $x' = (\vec{x}^T \cdot \vec{u}) \vec{u}$

Alternative criterion:
 Minimum squared distance
 $\|x - x'\|^2$

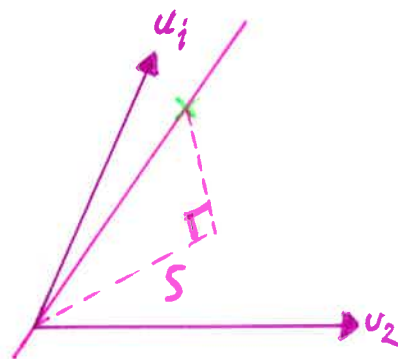
Subspace Method

Project a class onto a subspace of the dataset. then classify new points based on how far they are projected on the subspace.

$$S = \sum_{i=1}^p (x_i \cdot u_i)^2$$

\uparrow dimension of subspace
 \uparrow basis

large projection \Rightarrow more similar



Fisher's Criterion

The eigenvector for the greatest eigenvalue of $S_w^{-1} S_B$ gives A that maximizes fisher's criterion.

Ensemble Learning

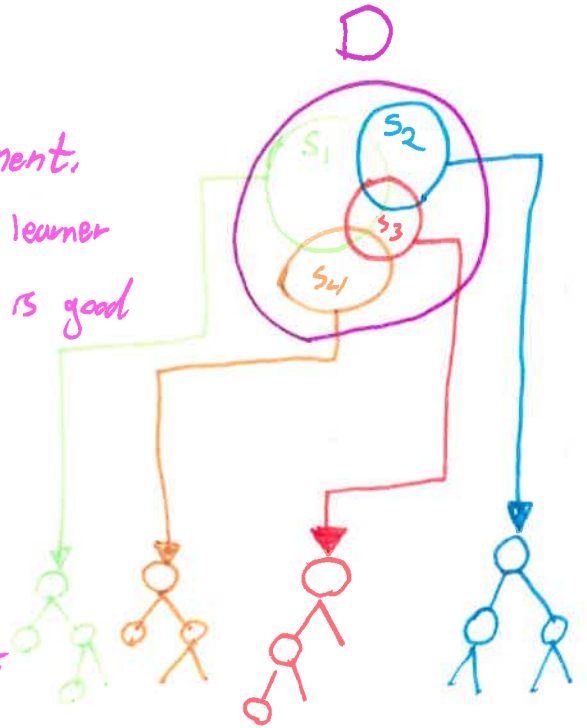
Combine multiple weak learners.
 Allows for each classifier to have higher variance.
 Assumes weak learners are independent

Diversity of opinion
 Independence
 Decentralization
 Aggregation.

Bootstrap Aggregating

Bagging

Train on dataset by sampling with replacement.
 Variance: High - Overtrain each individual learner
 Bias: Low - On average classification is good



Random Forest

Bagging + feature selection randomization

Instead of selecting best feature to split at each node, it is random so that trees are built different from each other.

Boosting

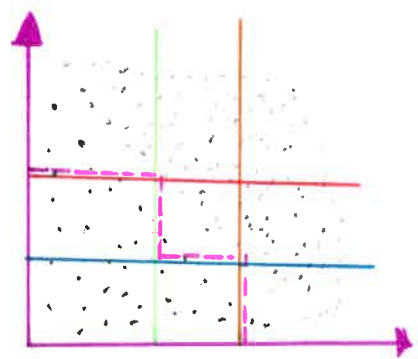
Treat misclassified samples as more important for each additional classifier

Variance: low - each model is very simple

Bias: High - Each learner is bad but if you combine them they are strong

Voting:
$$f(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

 (where α_t is the weight of each classifier and $h_t(x)$ is the individual weak learner)



Four weak learners combine to one stronger one.

Adaboost:

Lab 3

$$\alpha_t = \log_e \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

normalize to 1

$$W_i = w_i \cdot \exp(-\alpha_t y_i h_t(x_i))$$

Dropout

Randomly kill a node in an ANN and force other nodes to improve

