

Recap: classification

Statistical Methods in NLP 2

ISCL-BA-08

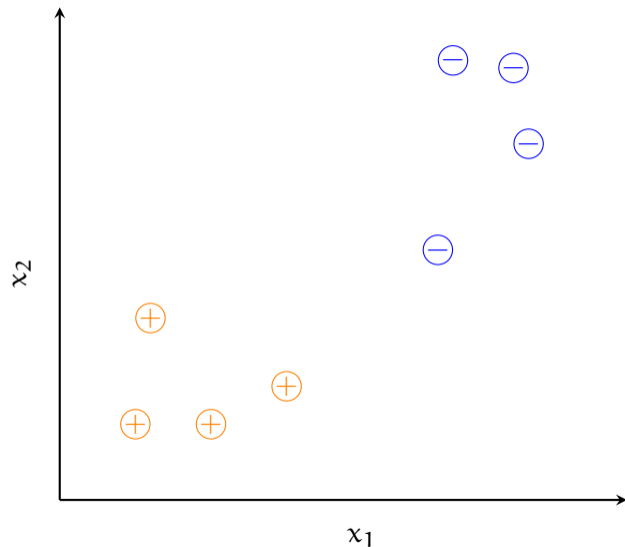
Çağrı Çöltekin
ccoltekin@sfs.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft

Summer Semester 2026

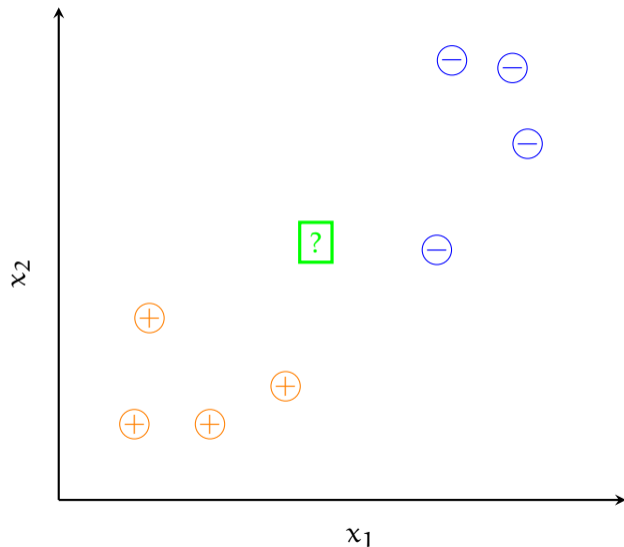
Classification

- Given a training set with (categorical) labels/outcome

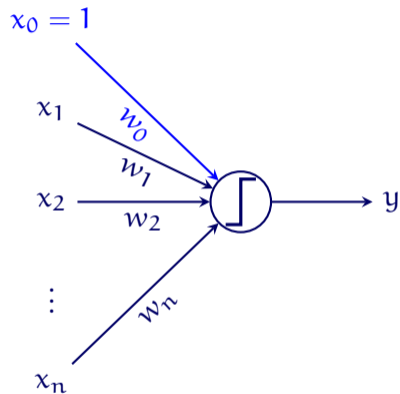


Classification

- Given a training set with (categorical) labels/outcome
- Train a model to predict future data points from the same distribution



The perceptron

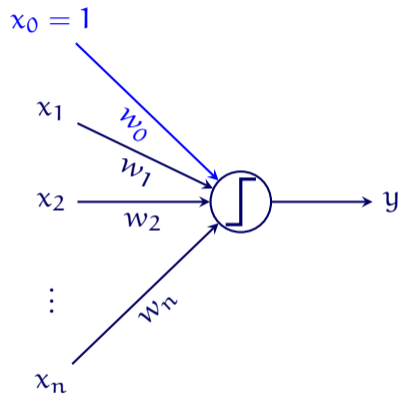


$$y = f \left(\sum_i^n w_i x_i \right)$$

where

$$f(x) = \begin{cases} +1 & \text{if } \sum_i^n w_i x_i > 0 \\ -1 & \text{otherwise} \end{cases}$$

The perceptron



$$y = f\left(\sum_i^n w_i x_i\right)$$

where

$$f(x) = \begin{cases} +1 & \text{if } \sum_i^n w_i x_i > 0 \\ -1 & \text{otherwise} \end{cases}$$

Note the similarity with linear regression: the separation between the (predicted) classes is a linear line/surface.

Learning with perceptron

- We do not update the parameters if classification is correct
- For misclassified examples, we try to minimize

$$E(\mathbf{w}) = - \sum_i \mathbf{w} \mathbf{x}_i y_i$$

where i ranges over all misclassified examples

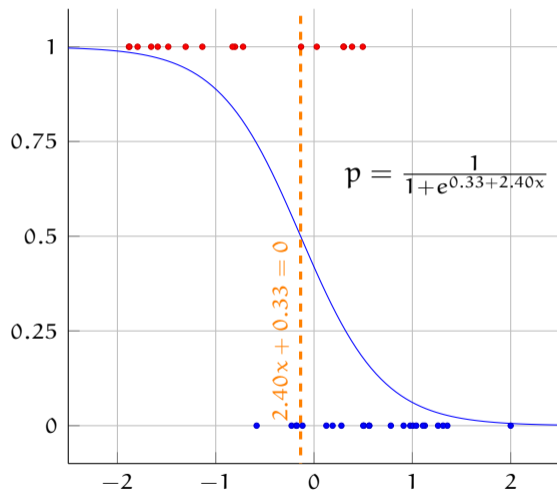
- Perceptron algorithm updates the weights for misclassified examples

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E(\mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{x}_i y_i$$

- Perceptron algorithm converges if the classes are *linearly separable*

Logistic regression

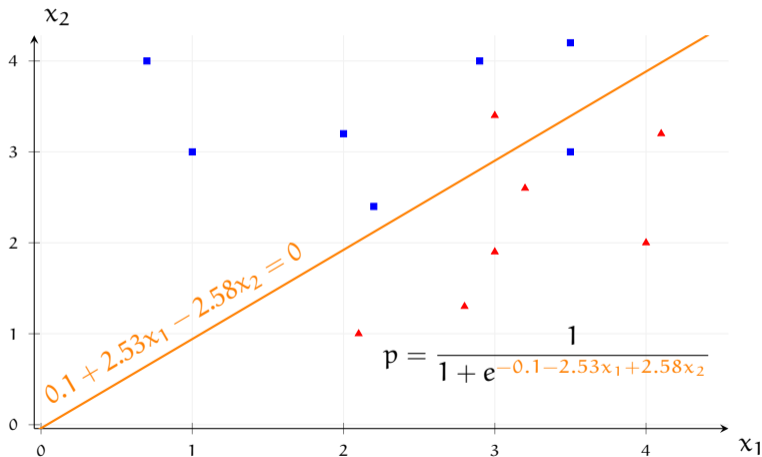


- Logistic regression is probabilistic, we estimate $p = P(y = 1|x)$
- Note that typical regression would penalize correctly classified examples
- Instead we fit a regression model for $\text{logit}(p) = \mathbf{w}^T \mathbf{x} + b$
- The probability estimate is the inverse of *logit*, the *logistic* (sigmoid) function

$$p = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

Another example with two predictors

Probability assignments are non-linear, but the discriminant function is linear



Fitting a logistic regression model

$$p = \frac{1}{1 + e^{-\mathbf{w}x}}$$

The likelihood of the training set is,

$$\mathcal{L}(\mathbf{w}) = \prod_i p^{y_i} (1 - p)^{1 - y_i}$$

In practice, we maximize log likelihood, or minimize '− log likelihood':

$$-\log \mathcal{L}(\mathbf{w}) = - \sum_i y_i \log p + (1 - y_i) \log(1 - p)$$

- The derivative/gradient is easy (a good exercise)
- There is no analytic solution for $\nabla - \log \mathcal{L}(\mathbf{w}) = 0$
- But the loss function is convex: we can find the global minimum with gradient descent (most of the time)

Naive Bayes

- *Naive Bayes* is another probabilistic classification method
- In any classification task, our aim is to find

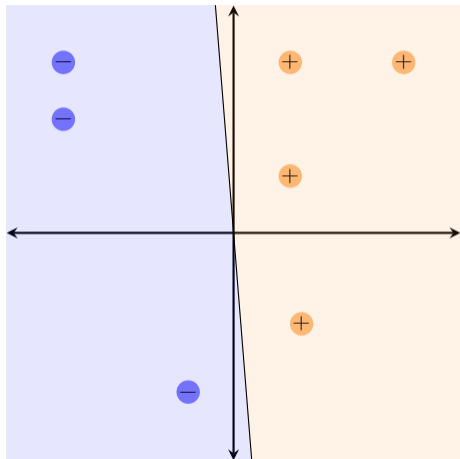
$$\hat{y} = \arg \max_y P(y | \mathbf{x})$$

- Sometimes (with some simplifying assumptions) it is easier to predict $P(\mathbf{x} | y)$
- Then we use Bayes' formula to invert the conditional probability

$$\hat{y} = \arg \max_y \frac{P(\mathbf{x} | y)P(y)}{P(\mathbf{x})} = \arg \max_y P(\mathbf{x} | y)P(y)$$

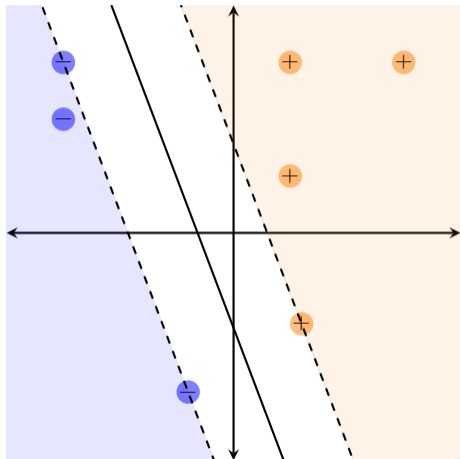
- $P(\mathbf{x} | y)$ and $P(y)$ are generally estimated using MLE (with *smoothing*)

Maximum-margin methods (e.g., SVMs)



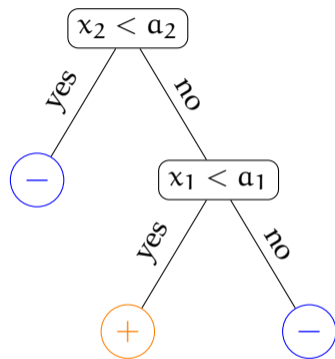
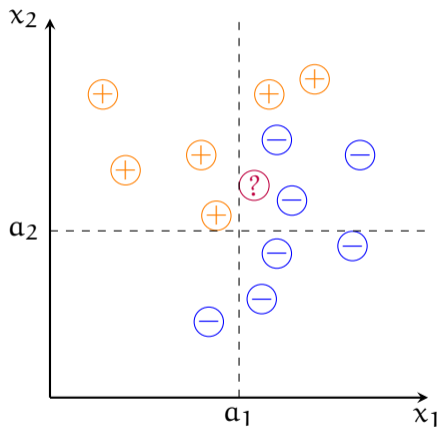
- In perceptron, we stop whenever we found a linear discriminator

Maximum-margin methods (e.g., SVMs)

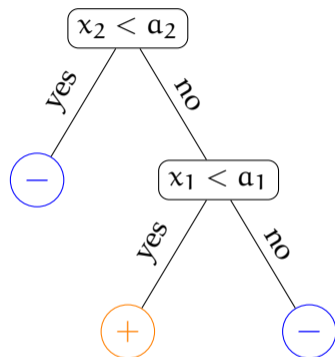
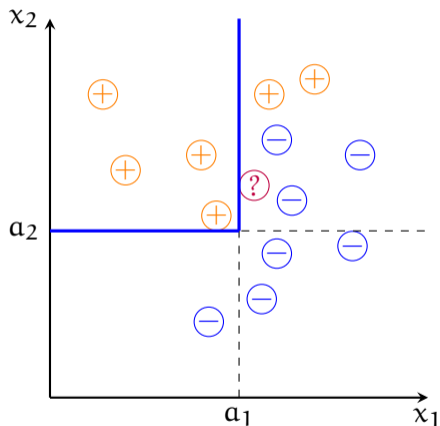


- In perceptron, we stop whenever we found a linear discriminator
- Maximum-margin classifiers seek a discriminator that maximizes the margin
- SVMs have other interesting properties, and they have been one of the best 'out-of-the-box' classifiers for many problems

Decision trees

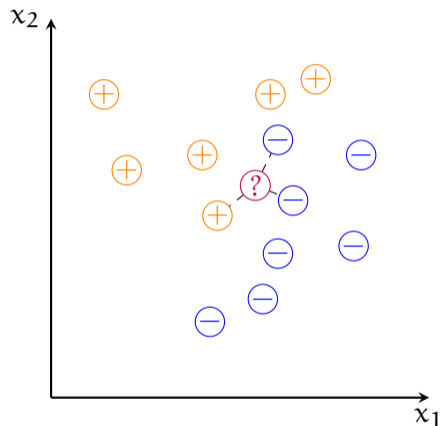


Decision trees



- Note that the decision boundary is non-linear

Instance/memory based methods



- No training: just memorize the instances
- During test time, decide based on the k nearest neighbors
- Like decision trees, **kNN** is non-linear
- It can also be used for regression

More than two classes

- Some algorithms can naturally be extended to handle multiple class labels
- Any binary classifier can be turned into a k-way classifier by

OvR **one-vs-rest** or **one-vs-all**

- train k classifiers: each learns to discriminate one of the classes from the others
- at prediction time the classifier with the highest confidence wins
- needs a confidence score from the base classifiers

OvO **one-vs-one**

- train $\frac{k(k-1)}{2}$ classifiers: each learns to discriminate a pair of classes
- decision is made by (weighted) majority vote
- works without need for confidence scores, but needs more classifiers

Measuring success in classification

Accuracy, precision, recall, F-score

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

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F}_1\text{-score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

		predicted	
		positive	negative
true value	pos.	TP	FN
	neg.	FP	TN

Multi-class evaluation

- For multi-class problems, it is common to report average precision/recall/f-score
- For C classes, averaging can be done two ways:

$$\text{precision}_M = \frac{\sum_i^C \frac{TP_i}{TP_i + FP_i}}{C} \quad \text{recall}_M = \frac{\sum_i^C \frac{TP_i}{TP_i + FN_i}}{C}$$

$$\text{precision}_\mu = \frac{\sum_i^C TP_i}{\sum_i^C TP_i + FP_i} \quad \text{recall}_\mu = \frac{\sum_i^C TP_i}{\sum_i^C TP_i + FN_i}$$

($M = \text{macro}$, $\mu = \text{micro}$)

- The averaging can also be useful for binary classification, if there is no natural positive class

Confusion matrix

- A confusion matrix is often useful for multi-class classification tasks

		predicted		
		negative	neutral	positive
true value	negative	10	2	0
	neutral	3	12	7
	positive	4	8	7

- Are the classes balanced?
- What is the accuracy?
- What is per-class, and averaged precision/recall?

Overfitting & Underfitting

We want our models to *generalize*, perform well on unseen data.

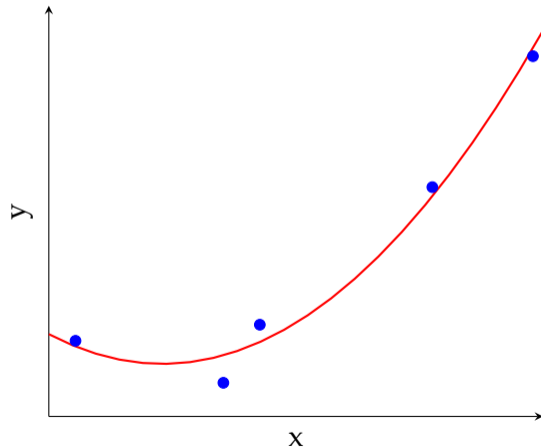
- *Overfitting* occurs when the model learns the idiosyncrasies of the training data
- *Underfitting* occurs when the model is not flexible enough for solving the problem at hand

We want simpler models, but not too simple for the task at hand.

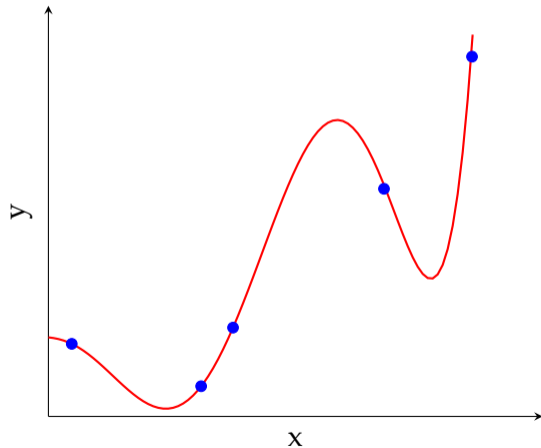
Causes of overfitting

(1) Model complexity

Quadratic



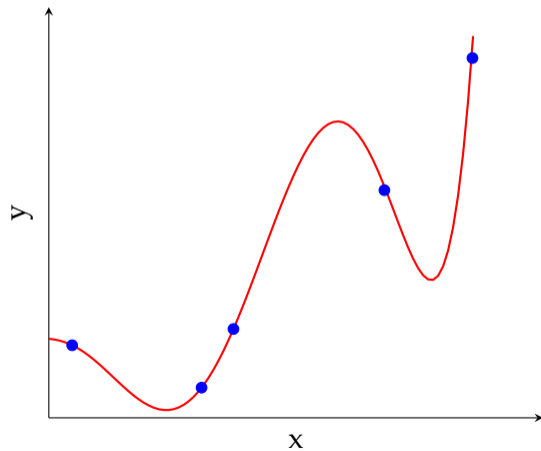
6th degree



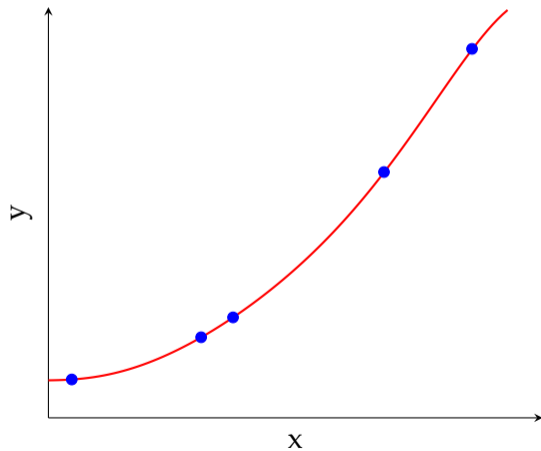
Causes of overfitting

(1) Noise

With noise



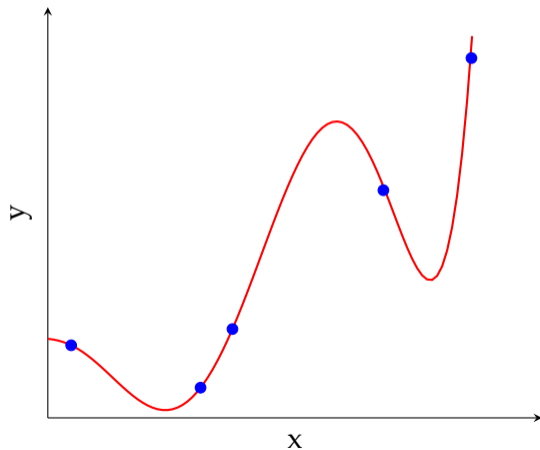
No noise



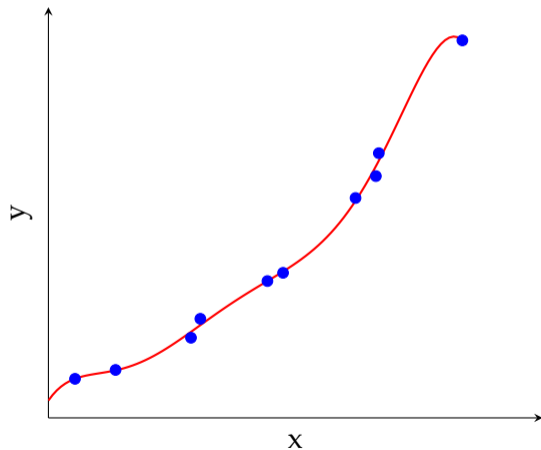
Causes of overfitting

(1) Lack of training data

5 training instances



10 training instances



Preventing overfitting

- Add more training data
- Increase data quality
- Lower model complexity
 - Simpler models
 - Regularization
 - For iterative methods, stop early
- In most cases, selecting the best model requires *hyperparameter* tuning / optimization.
- Always check the model performance on a held-out (development / validation) set.
- K-fold cross-validation is a method for efficient use of available labeled data

Bias and variance

Bias of an estimate is the difference between the value being estimated, and the expected value of the estimate

$$B(\hat{\mathbf{w}}) = E[\hat{\mathbf{w}}] - \mathbf{w}$$

- An *unbiased* estimator has 0 bias

Variance of an estimate is, simply its variance, the value of the squared deviations from the mean estimate

$$\text{var}(\hat{\mathbf{w}}) = E \left[(\hat{\mathbf{w}} - E[\hat{\mathbf{w}}])^2 \right]$$

\mathbf{w} is the parameter (vector) that defines the model

Bias–variance relationship is a trade-off:
models with low bias result in high variance.

ML evaluation in general

The first principle is that you must not fool yourself and you are the easiest person to fool.
– Richard P. Feynman

- We want models with low bias and low variance, but this is a trade-off
- Estimators/models with high variance (hence, low/no bias) are likely to *overfit*
- Estimators/models with low variance (hence, high bias) are likely to *underfit*
- Evaluating ML system requires special care:
 - Tuning your system on a development set
 - Cross-validation allows efficient use of labeled data during tuning
 - A test set is often used when comparing results obtained by different models
- In most cases, the scores are not meaningful by themselves, we need to compare with *baselines*

Final remarks

- Most NLP problems we try to solve are classification problems
- We reviewed some of the ‘traditional’ classification methods
- Understanding them will help understanding more ‘modern’ methods
- The models we review can serve as baselines for more complex models, and sometimes their performance may surprise you

Final remarks

- Most NLP problems we try to solve are classification problems
- We reviewed some of the ‘traditional’ classification methods
- Understanding them will help understanding more ‘modern’ methods
- The models we review can serve as baselines for more complex models, and sometimes their performance may surprise you

Next:

Mon/Wed Gradient descent, reading: Jurafsky and Martin (2026, Section 4.6)


Fri Lab: numpy tutorial

Some sources of information


- Any ML textbook covers most of the methods reviewed (and more), here are a few: James et al. (2024), Bishop (2006), and MacKay (2003)

 Bishop, Christopher M. (2006). *Pattern Recognition and Machine Learning*. Springer. ISBN: 978-0387-31073-2.

 James, Gareth, Daniela Witten, Trevor Hastie, Robert Tibshirani, and Jonathan Taylor (2024). *An introduction to statistical learning*. Springer. ISBN: 9783031391897. URL: <https://www.statlearning.com/>.

 Jurafsky, Daniel and James H. Martin (2026). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd. Online manuscript released January 6, 2026. URL: <https://web.stanford.edu/~jurafsky/slp3/>.

Some sources of information (cont.)

 MacKay, David J. C. (2003). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press. ISBN: 978-05-2164-298-9. URL: <http://www.inference.phy.cam.ac.uk/itprnn/book.html>.