

Machine Learning

Linear Models for Classification II

Manfred Jaeger

Aalborg University

Probabilistic Models

Classify data instance \mathbf{x} to belong to the class k , for which probability

$$P(Y = k \mid \mathbf{x})$$

is maximal.

Classify data instance \mathbf{x} to belong to the class k , for which probability

$$P(Y = k \mid \mathbf{x})$$

is maximal.

Generative Approach

- ▶ Define a **joint distribution** for the features $\mathbf{X} = X_1, \dots, X_D$ and the class variable Y with values $1, \dots, K$. Usually factored as $P(Y, \mathbf{X}) = P(Y)P(\mathbf{X} \mid Y)$.
- ▶ Compute

$$P(Y = k \mid \mathbf{X} = \mathbf{x}) = \frac{P(Y = k, \mathbf{X} = \mathbf{x})}{P(\mathbf{X} = \mathbf{x})} = P(\mathbf{X} = \mathbf{x} \mid Y = k) \frac{P(Y = k)}{P(\mathbf{X} = \mathbf{x})}$$

Example: Naive Bayes, LDA

Classify data instance \mathbf{x} to belong to the class k , for which probability

$$P(Y = k \mid \mathbf{x})$$

is maximal.

Generative Approach

- ▶ Define a **joint distribution** for the features $\mathbf{X} = X_1, \dots, X_D$ and the class variable Y with values $1, \dots, K$. Usually factored as $P(Y, \mathbf{X}) = P(Y)P(\mathbf{X} \mid Y)$.
- ▶ Compute

$$P(Y = k \mid \mathbf{X} = \mathbf{x}) = \frac{P(Y = k, \mathbf{X} = \mathbf{x})}{P(\mathbf{X} = \mathbf{x})} = P(\mathbf{X} = \mathbf{x} \mid Y = k) \frac{P(Y = k)}{P(\mathbf{X} = \mathbf{x})}$$

Example: Naive Bayes, LDA

Discriminative Approach

Directly learn the conditional distribution

$$P(Y \mid \mathbf{X})$$

Example: Logistic Regression, (Neural Networks)

Operations supported:

	Generative	Discriminative
Predicting Y given \mathbf{X}	yes	yes
Predicting any other variable X_i	yes (compute arbitrary conditional probabilities $P(X_i \dots)$)	no
Generating data points (\mathbf{x}, y)	yes	no

- ▶ Operations that are supported (in principle) by generative models are not always supported in a computationally efficient manner.
- ▶ For a fixed prediction task, discriminative models can be more efficient to learn (w.r.t the required data), and they can be more interpretable.

Notation: Upper case X, \mathbf{X} for (tuples of) *variables*. Lower case x, \mathbf{x} for (tuples of) specific values that the variables can have. **Example:** $X = \text{Temperature}$, $x = 23.5^\circ\text{C}$.

In all cases, the main criterion for learning probabilistic models is maximizing the likelihood:

Generative

$$\prod_{n=1}^N P(\mathbf{x}_n, y_n)$$

Discriminative

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

Always required: a **parametric model** for the (conditional) distributions, i.e. specific functional forms

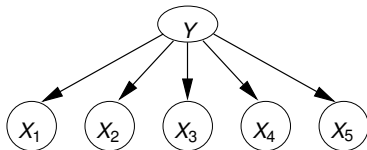
$$P(Y \mid \mathbf{w}), P(\mathbf{X} \mid Y, \mathbf{w}), P(Y \mid \mathbf{X}, \mathbf{w})$$

defined in terms of parameter vectors \mathbf{w} .

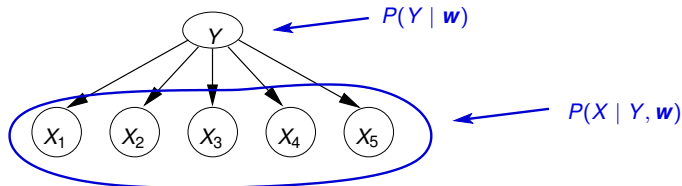
Model $P(Y \mid \mathbf{w})$ usually **multinomial** (“**multinoulli**”), i.e.

$$\mathbf{w} = (P(Y = 1), P(Y = 2), \dots, P(Y = K - 1))$$

The naive Bayes model is a probabilistic generative model:



The naive Bayes model is a probabilistic generative model:



$P(Y | \mathbf{w})$: Conditional probability table for node Y ; \mathbf{w} : entries in the conditional probability table

$P(X_i | Y, \mathbf{w})$:
▶ if X_i categorical (discrete): conditional probability tables with entries defined by \mathbf{w}
▶ if X_i continuous: e.g.: every distribution $P(X_i | Y, \mathbf{w})$ is a Gaussian distribution with parameters (mean, variance) defined by \mathbf{w} .

➡ a simple Gaussian mixture model

Gaussian Mixtures

The Gaussian mixture model for a D -dimensional feature space, without the naive Bayes independence assumptions:

- ▶ $P(Y)$ multinomial
- ▶ $P(\mathbf{X} | Y)$ Gaussian, i.e. for all $k = 1, \dots, K$:

$$P(\mathbf{X} = \mathbf{x} | Y = k) = \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}$$

defined by parameters $\boldsymbol{\mu}_k$ (mean vectors) and Σ_k (co-variance matrices).

Example

- ▶ $Y = \text{Region} \in \{\text{Nordjylland}, \dots, \text{Hovedstaden}\}$ ($K = 5$).
- ▶ $\mathbf{X} = (X_1, X_2)$ with X_1 : annual income, X_2 : annual housing expenditure (mortgage, rent).
- ▶ $\boldsymbol{\mu}_k = (\mu_{1,k}, \mu_{2,k})$: average income and housing costs in region k .
- ▶ $\Sigma_k = \begin{pmatrix} \sigma_{1,1,k} & \sigma_{1,2,k} \\ \sigma_{2,1,k} & \sigma_{2,2,k} \end{pmatrix}$

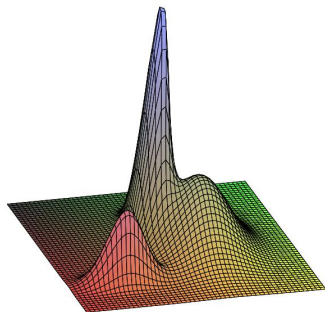
➡ Presumably: X_1 and X_2 are positively correlated, i.e., $\sigma_{1,2,k} (= \sigma_{2,1,k}) > 0$.

Example with $D = 2$ and diagonal covariance matrices (= naive Bayes!):

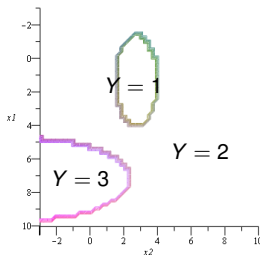
$$P(Y = 1) = 0.3, \quad P(Y = 2) = 0.6, \quad P(Y = 3) = 0.1$$

$$\mu_1 = (2, 3), \quad \mu_2 = (4, 5), \quad \mu_3 = (7, 0)$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 0.25 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 4 & 0 \\ 0 & 2.25 \end{pmatrix}, \quad \Sigma_3 = \begin{pmatrix} 0.25 & 0 \\ 0 & 1 \end{pmatrix}$$



Distribution of \mathbf{X}



Decision boundaries

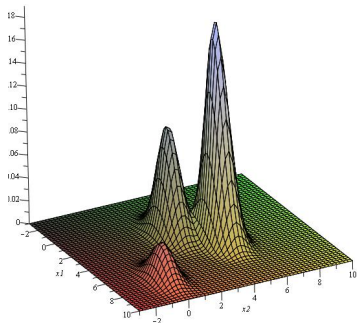
LDA: special case where all covariance matrices are identical. [Warning: not to be confused with *Fisher's linear discriminant* – Bishop Sec. 4.2.1 describes LDA, but does not use this name!]

Example:

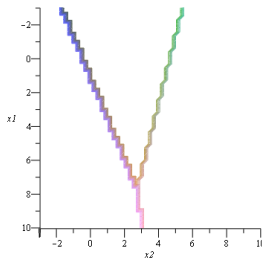
$$P(Y = 1) = 0.3, P(Y = 2) = 0.6, P(Y = 3) = 0.1$$

$$\mu_1 = (2, 3), \mu_2 = (4, 5), \mu_3 = (7, 0)$$

$$\Sigma_1 = \Sigma_2 = \Sigma_3 = \begin{pmatrix} 1 & 0 \\ 0 & 0.25 \end{pmatrix}$$

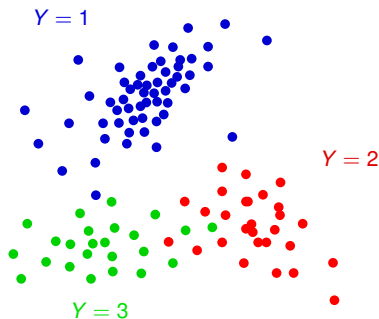


Distribution of \mathbf{X}



Decision boundaries linear!

Observed data:



Total of N data points (\mathbf{x}_n, y_n) ;
 N_k data points belonging to class k

Maximum likelihood parameters:

- ▶ $P(Y = k) := \frac{N_k}{N}$
- ▶ $\mu_k := \frac{1}{N_k} \sum_{n: y_n = k} \mathbf{x}_n$
- ▶ $\Sigma_k[i, j] := \frac{1}{N_k} \sum_{n: y_n = k} (\mathbf{x}_n[i] - \mu_k[i])(\mathbf{x}_n[j] - \mu_k[j])$ (unrestricted mixtures)
- ▶ $\Sigma[i, j] := \frac{1}{N} \sum_k \sum_{n: y_n = k} (\mathbf{x}_n[i] - \mu_k[i])(\mathbf{x}_n[j] - \mu_k[j])$ (LDA)

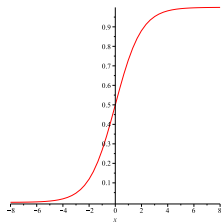
Logistic Regression

Assume continuous features \mathbf{X} , and two classes: $Y \in \{0, 1\}$.

Define

$$P(Y = 1 \mid \mathbf{X} = \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w} \cdot \mathbf{x})$$

where \mathbf{w} is a weight vector, and σ is the sigmoid function:



$$\sigma(x) = \frac{e^x}{1+e^x} = \frac{1}{1+e^{-x}}$$

Then

$$P(Y = 1 \mid \mathbf{X} = \mathbf{x}, \mathbf{w}) \geq 0.5 \Leftrightarrow \mathbf{w} \cdot \mathbf{x} \geq 0$$

\rightsquigarrow linear decision boundary!

➡ Not a generative model: we cannot use model as a random generator of (\mathbf{x}, y) data.

Categorical Features

- ▶ Categorical feature X with values $\{x_1, \dots, x_h\}$ is coded by h “numerical” features via one-hot-encoding.

Categorical Features

- ▶ Categorical feature X with values $\{x_1, \dots, x_h\}$ is coded by h “numerical” features via one-hot-encoding.

More than 2 classes

Alternative view of logistic regression formula:

$$\frac{P(Y = 1 \mid \mathbf{X} = \mathbf{x}, \mathbf{w})}{P(Y = 0 \mid \mathbf{X} = \mathbf{x}, \mathbf{w})} = e^{\mathbf{w} \cdot \mathbf{x}} \quad \text{or} \quad \log\left(\frac{P(Y = 1 \mid \mathbf{X} = \mathbf{x}, \mathbf{w})}{P(Y = 0 \mid \mathbf{X} = \mathbf{x}, \mathbf{w})}\right) = \mathbf{w} \cdot \mathbf{x}$$

▶ The linear function $\mathbf{w} \cdot \mathbf{x}$ defines the *log-odds* of class 1 against class 0.

For multi-class: $Y \in \{1, \dots, K\}$:

- ▶ Pick a *reference class*, e.g., $Y = K$
- ▶ Learn $K - 1$ models (= parameters \mathbf{w}_k) for the (log) odds
 $odds(k, K, \mathbf{x}) = P(Y = k \mid \mathbf{X} = \mathbf{x}, \mathbf{w}_k) / P(Y = K \mid \mathbf{X} = \mathbf{x}, \mathbf{w}_k)$ ($k = 1, \dots, K - 1$)
- ▶ Classify point \mathbf{x} as the k for which $odds(k, K, \mathbf{x})$ is maximal, if this is > 1 , and as K otherwise.

▶ K discriminant functions: $odds(k, K, \mathbf{x})$ ($k = 1, \dots, K - 1$) and $odds(K, K, \mathbf{x}) \equiv 1$.

Maximize likelihood

$$\prod_{n:y_n=1} P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}) \prod_{n:y_n=0} (1 - P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}))$$

Taking the log: maximize

$$\sum_{n:y_n=1} \log P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}) + \sum_{n:y_n=0} \log (1 - P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}))$$

Rewritten, using the 0, 1 encoding of the two classes:

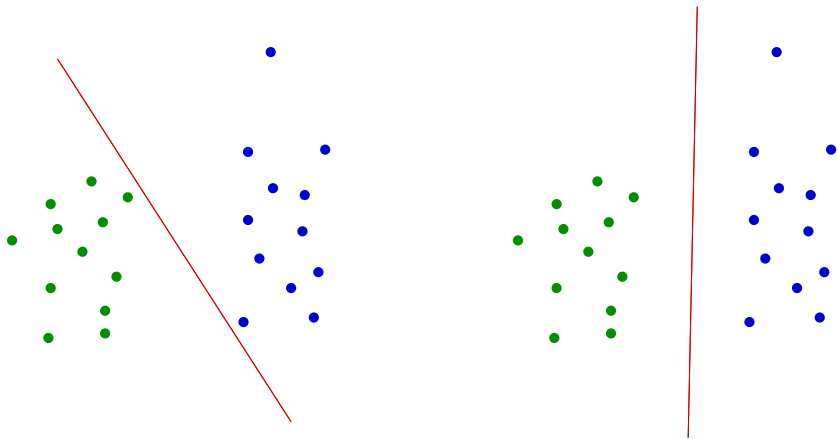
$$\sum_{n=1}^N y_n \log P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}) + (1 - y_n) \log (1 - P(Y = 1 \mid \mathbf{X} = \mathbf{x}_n, \mathbf{w}))$$

- Set derivative to zero; find root of derivative using numerical methods, often Newton-Raphson

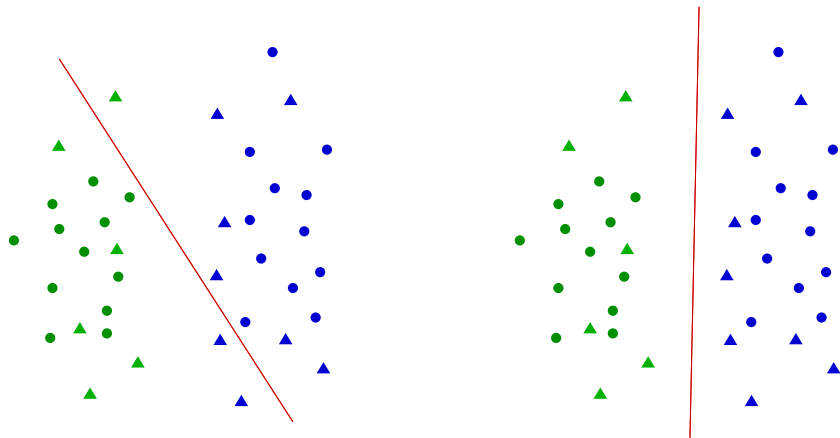
	Least Squares	Perceptron	LDA	Logistic Regression
Criterion	Approximate one-hot class encoding	Minimize error function	Maximize Likelihood	Maximize Likelihood
Multi-class	Yes	No, but extendible	yes	No, but extendible
Learning	Solving matrix equation	Iterative optimization	One-step optimization	Iterative optimization
Finds separating decision boundary if classes linearly separable ($K = 2$):	not always	yes	not always	not always
Works for not linearly separable data	yes	poorly (may not converge)	yes	yes

Maximum Margin Hyperplanes

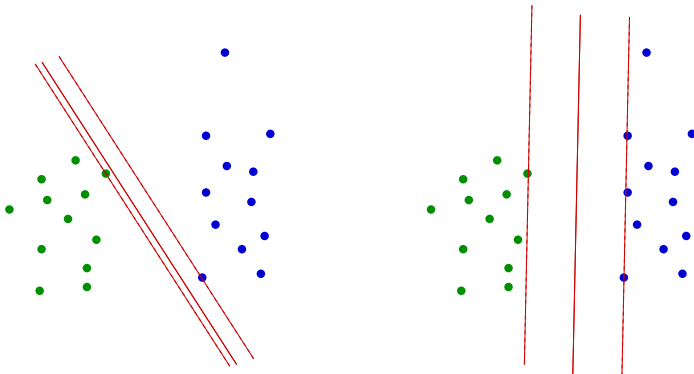
Two separating hyperplanes:



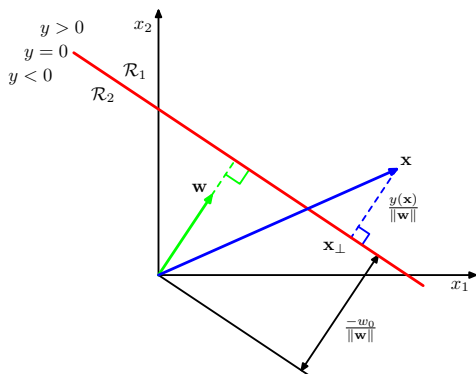
Two separating hyperplanes:



Hyperplane on the right expected to have lower error on new test data (triangles)



- ▶ Margin of a separating hyperplane: minimum distance of any datapoint to the hyperplane
- ▶ **Maximum-margin hyperplane**: hyperplane with maximum margin.
- ▶ For Max-margin hyperplane:
 - ▶ minimum distance of point from green class to hyperplane = minimum distance of point from blue class to hyperplane
 - ▶ there are datapoint in both classes whose distance to the hyperplane equals the margin. These are the **support vectors** of the hyperplane



$$y(\mathbf{x}) = w_0 + \mathbf{w} \cdot \mathbf{x}$$

Distance of \mathbf{x} to hyperplane $y = 0$:

$$\frac{|y(\mathbf{x})|}{\|\mathbf{w}\|}$$

(in Figure: $y(\mathbf{0}) = w_0 < 0$ and $y(\mathbf{x}) > 0$, so $y(\mathbf{x}) = |y(\mathbf{x})|$)

[Bishop, Fig. 4.1]

With class encoding $Y \in \{-1, 1\}$:

distance of datapoint \mathbf{x}_n with label y_n to hyperplane $y = 0$ (if datapoint lies on correct side of the hyperplane, i.e., $y_n = 1$ and $y(\mathbf{x}_n) > 0$, or $y_n = -1$ and $y(\mathbf{x}_n) < 0$): $\frac{y_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$

("negative distance" for points that lie on the wrong side of the hyperplane)

(Assume that data is linearly separable!)

Find hyperplane defined by \mathbf{w} , b as

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] \right\}$$

(b formerly known as w_0).

Simplification

The decision boundary (and distance of datapoints to decision boundary) is not changed by scaling \mathbf{w} , b with a common factor κ :

$$\mathbf{w} \mapsto \kappa \mathbf{w}, \quad b \mapsto \kappa b$$

\rightsquigarrow can calibrate \mathbf{w} , b so that for support vectors \mathbf{x}_n :

$$y_n(\mathbf{w} \cdot \mathbf{x}_n + b) = 1,$$

and for *all* datapoints \mathbf{x}_n :

$$y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1.$$

The original optimization problem

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] \right\}$$

then becomes: find

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|}$$

subject to the constraints

$$y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \quad (n = 1, \dots, N).$$

SVM Learning Objective

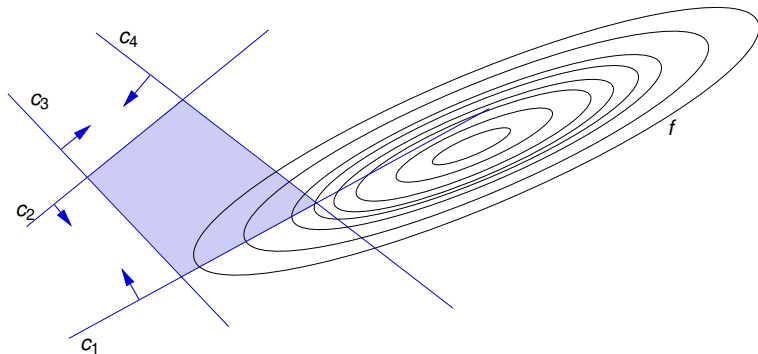
Replacing maximization of $1 / \|\mathbf{w}\|$ by equivalent minimization of $\|\mathbf{w}\|^2 / 2$ leads to the formulation: find

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

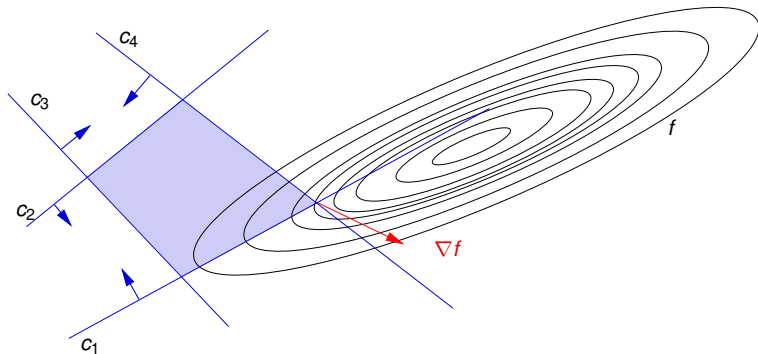
subject to the constraints

$$y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \quad (n = 1, \dots, N).$$

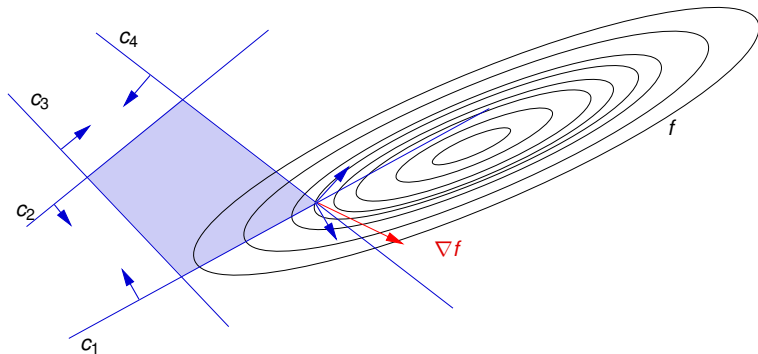
Constrained optimization problem: maximize $f(\mathbf{w})$ subject to *constraints* $c_i(\mathbf{w}) \geq 0$:



Constrained optimization problem: maximize $f(\mathbf{w})$ subject to *constraints* $c_i(\mathbf{w}) \geq 0$:



Constrained optimization problem: maximize $f(\mathbf{w})$ subject to *constraints* $c_i(\mathbf{w}) \geq 0$:



At the constrained optimum, we have the following relationship between the gradients of f and the gradients of the constraints:

$$\nabla f = \sum -\lambda_i \nabla c_i \text{ with } \lambda_i \geq 0$$

$$\nabla f = \sum -\lambda_i \nabla c_i \Leftrightarrow \nabla(f + \sum \lambda_i c_i) = 0$$

where $f + \sum \lambda_i c_i$ is the **Lagrange function** of the problem, and the λ_i are the **Lagrange multipliers**.

- ▶ The technique of Lagrange multipliers transforms the constrained optimization problem for \mathbf{w} into an unconstrained optimization problem for $\boldsymbol{\lambda}$ and \mathbf{w} .
- ▶ The solution for $\boldsymbol{\lambda}$ identifies the constraints c_i that are *active* at the solution: those with $\lambda_i > 0$.

The SVM learning problem is solved using the method of Lagrange multipliers.

- ▶ the solution identifies the **support vectors**: \mathbf{x}_i that lie on the margin (= those \mathbf{x}_i whose constraint is active at the optimal solution)
- ▶ the vector \mathbf{w} is a linear combination of support vectors:

$$\mathbf{w} = \sum_{i:\lambda_i>0} \lambda_i y_i \mathbf{x}_i$$

- ▶ a test instance \mathbf{z} is classified as

$$\text{sign}(\mathbf{w} \cdot \mathbf{z} + b) = \text{sign}\left(\sum_{i:\lambda_i>0} \lambda_i y_i \mathbf{x}_i \cdot \mathbf{z} + b\right)$$

- ▶ In the Lagrange optimization process: to determine the support vectors \mathbf{x}_i , the λ_i and b : the only operations required on data items is to compute dot products $\mathbf{x}_i \cdot \mathbf{x}_j$.
- ▶ for classification: only need to compute dot products $\mathbf{x}_i \cdot \mathbf{z}$

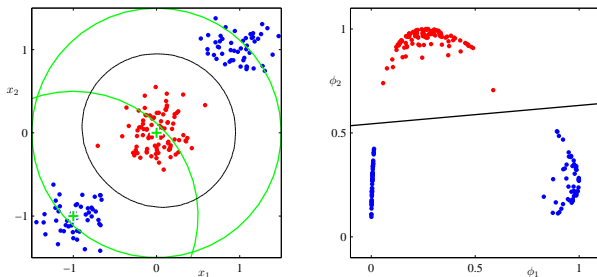
Data Transformations

Any mapping

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$$

defines a data transformation that transforms the original data instance \mathbf{x} into a transformed data instance $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_{D'}(\mathbf{x}))$.

- ▶ The components $\phi_i(\mathbf{x})$ are also called **features** or **basis functions**, and $\mathbb{R}^{D'}$ the **feature space** of ϕ .
- ▶ Often: $D' > D$, and ϕ_i typically non-linear.



[Bishop Fig. 4.12]

Example: Original data (left), transformed data in feature space (right), decision boundaries (black lines). Here: $\phi_1(x_1, x_2) = (x_1 + 1)^2 + (x_2 + 1)^2$ (?) and $\phi_2(x_1, x_2) = 1 - (x_1^2 + x_2^2)$ (?)

- All our linear models can be applied to the transformed data
- The linear decision boundary in feature space then corresponds to a non-linear decision boundary in the original space