

Machine Learning: The Instance-Based Perspective

Konstantin Tretyakov

<http://kt.era.ee>

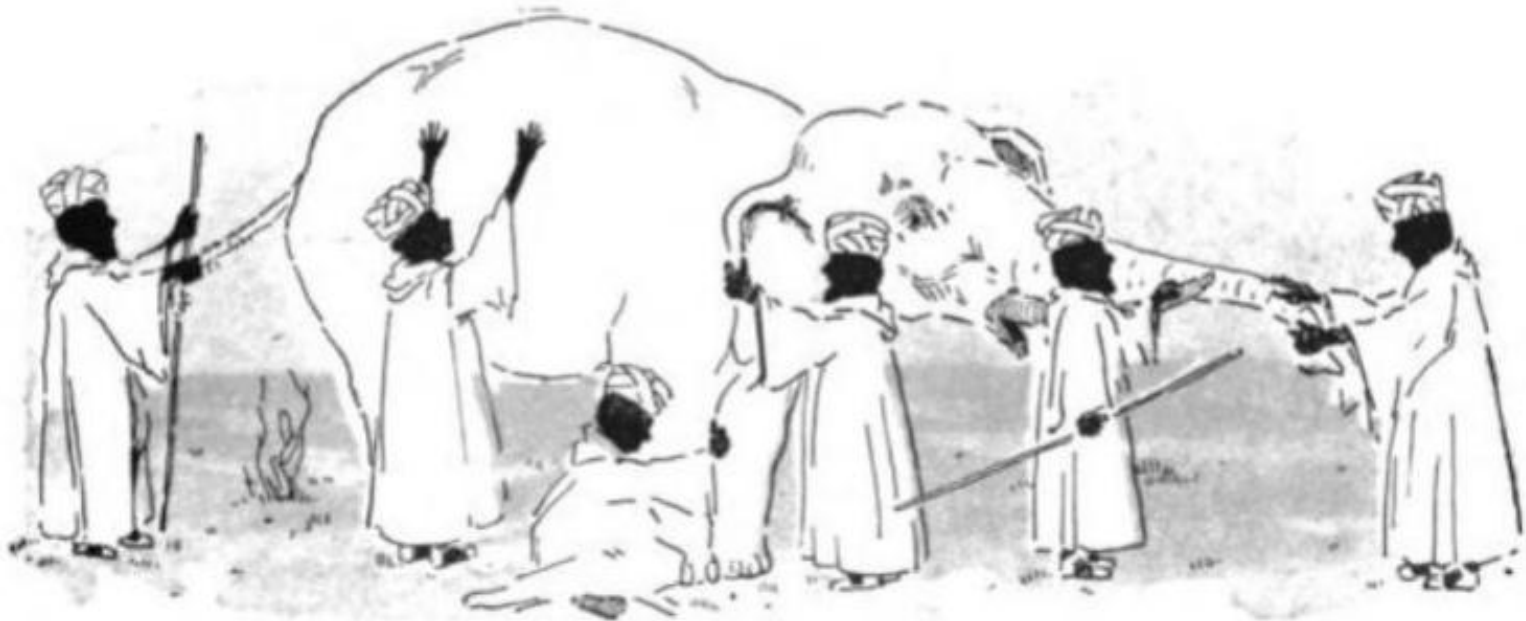
AACIMP 2012
August 2012

STACC

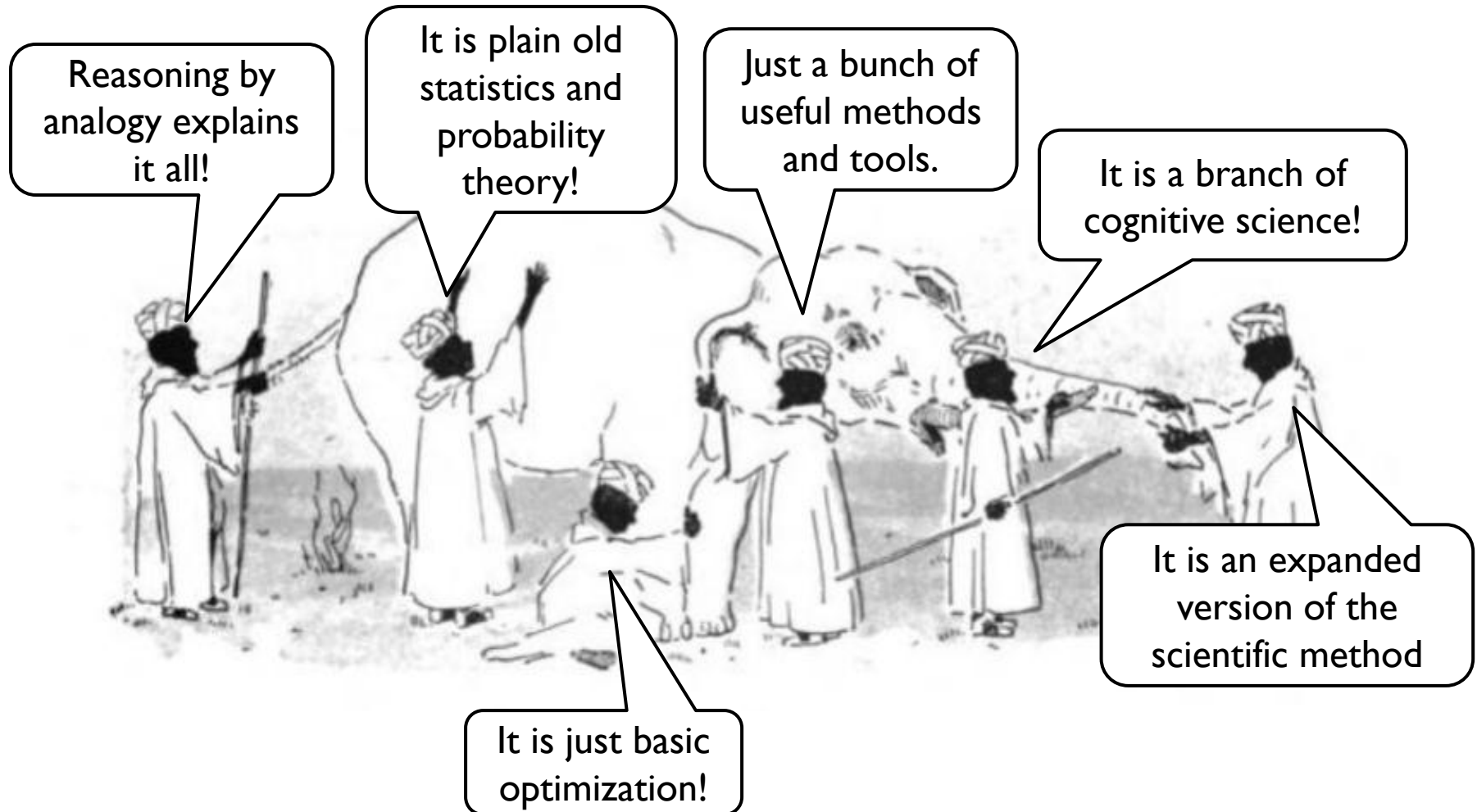
Software Technology and
Applications Competence Center



So far...



So far...



Today

Reasoning by
analogy explains
it all!

It is plain old
statistics and
probability
theory!

Just a bunch of
useful methods
and tools.

It is a branch of
cognitive science!



Linear classification & SVMs

Kernel Methods

an expanded
version of the
specific method

Next

Reasoning by
analogy explains
it all!

It is plain old
statistics and
probability
theory!

Just a bunch of
useful methods
and tools.

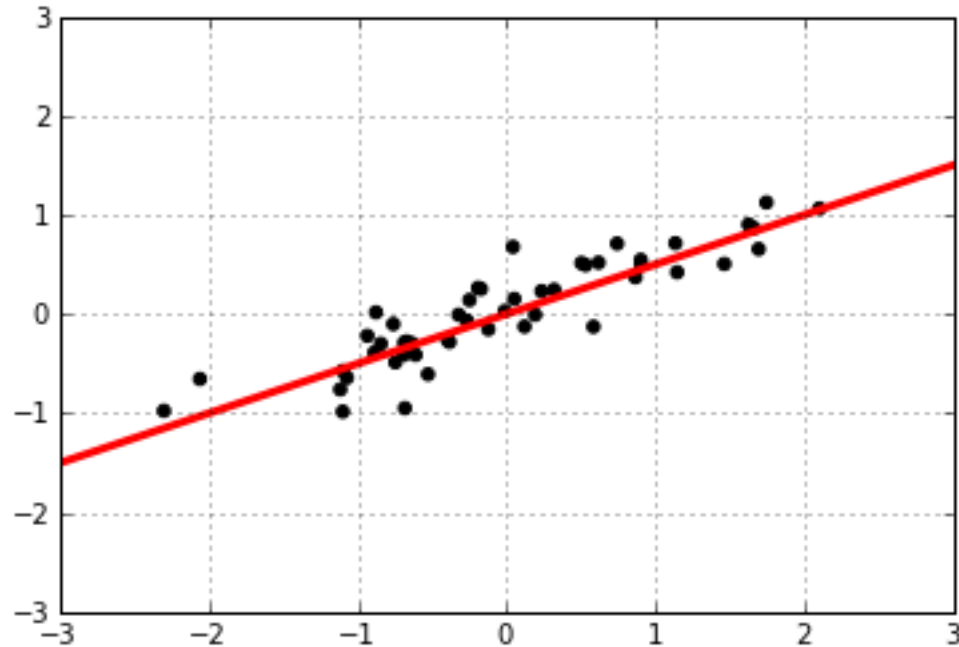
It is a branch of
cognitive science!

Linear classification & SVMs

Kernel Methods

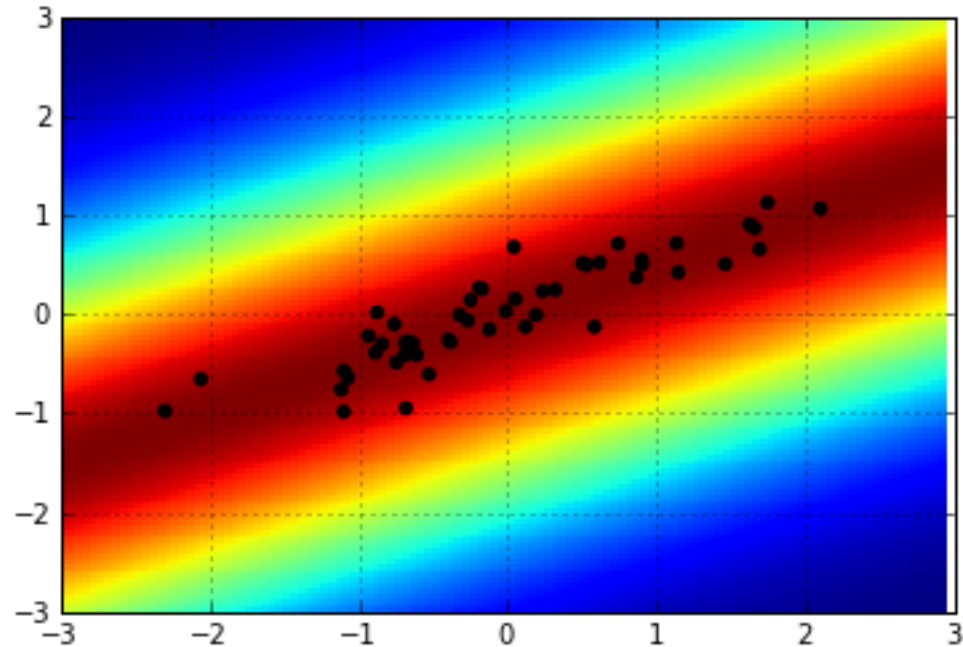
an expanded
version of the
specific method

Recall linear regression



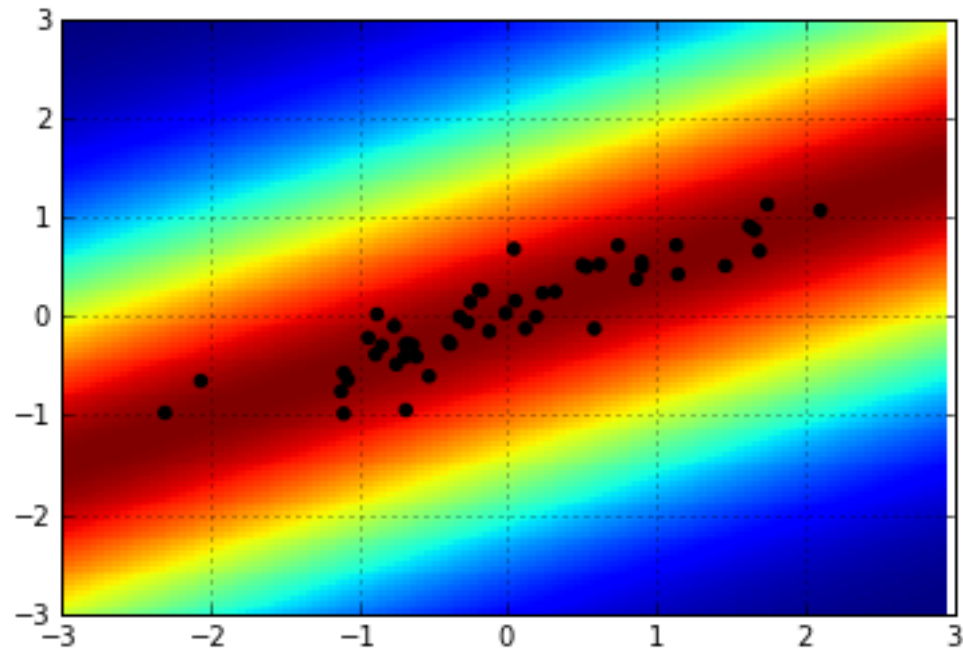
$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

Recall linear regression



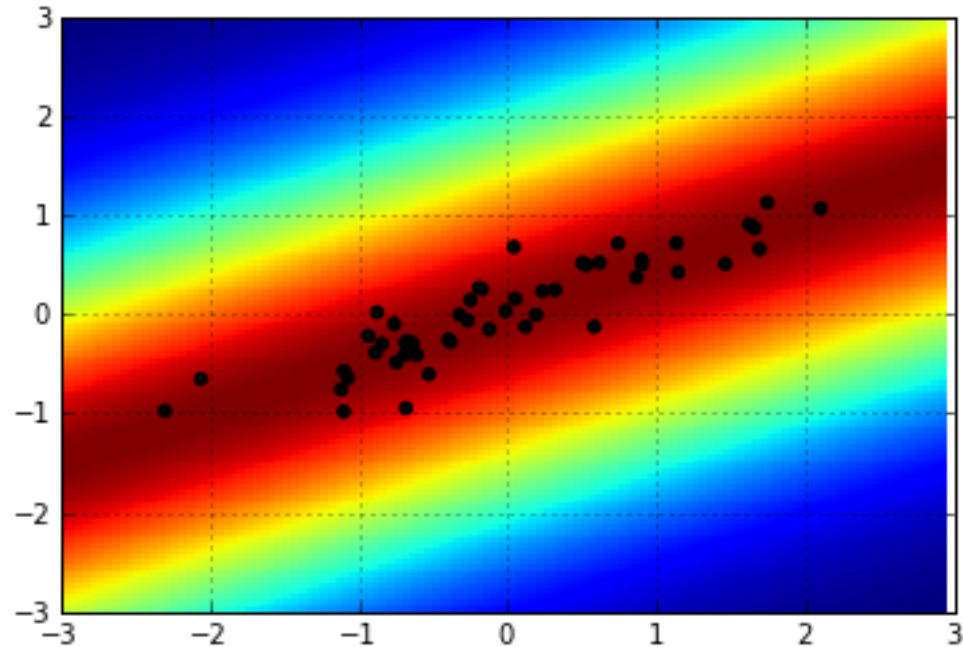
$$\hat{y} = \mathbf{w}^T \mathbf{x}$$
$$E(y|\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Recall linear regression



$$\hat{y} = \mathbf{w}^T \mathbf{x}$$
$$E(y|\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$
$$Y \sim \mathbf{w}^T \mathbf{x} + N(0, \sigma^2)$$

Recall linear regression



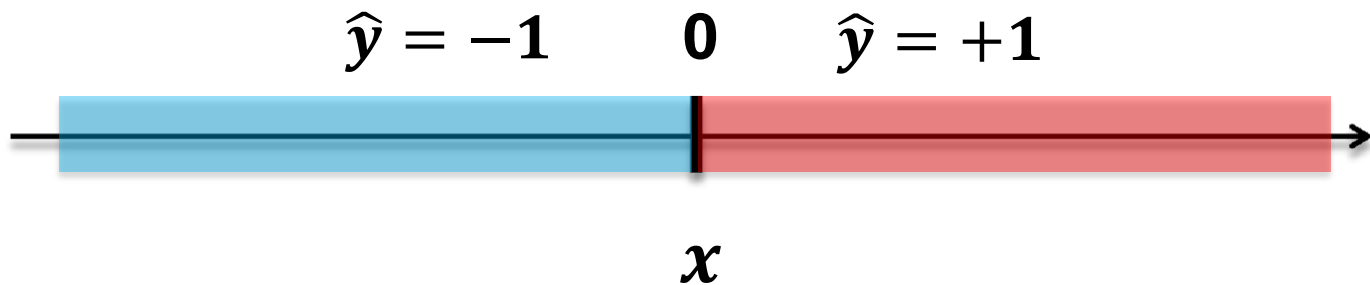
$$\hat{y} = \mathbf{w}^T \mathbf{x}$$
$$E(y|x) = \mathbf{w}^T \mathbf{x}$$
$$Y \sim N(\mathbf{w}^T \mathbf{x}, \sigma^2)$$

Linear classification

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

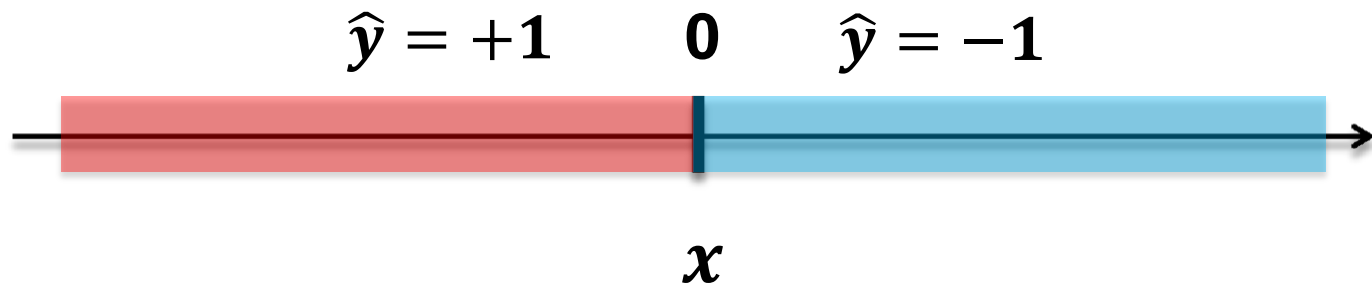
Linear classification

$$\hat{y} = \text{sign}(5 \cdot x)$$



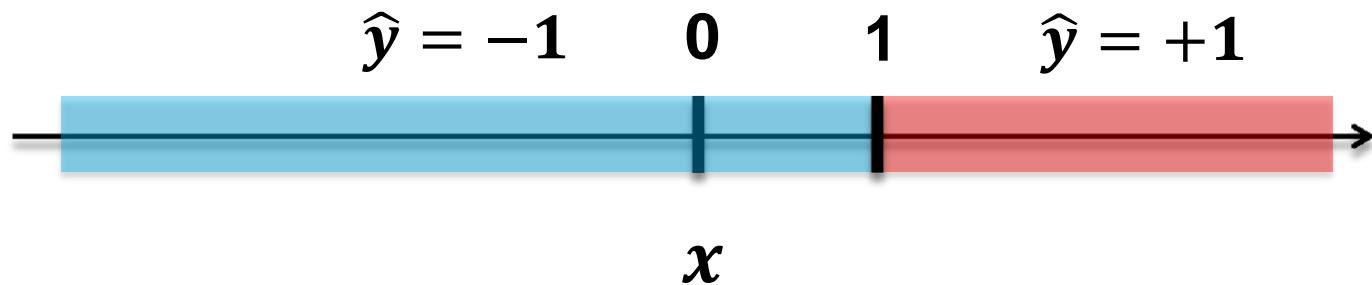
Linear classification

$$\hat{y} = \text{sign}(-5 \cdot x)$$



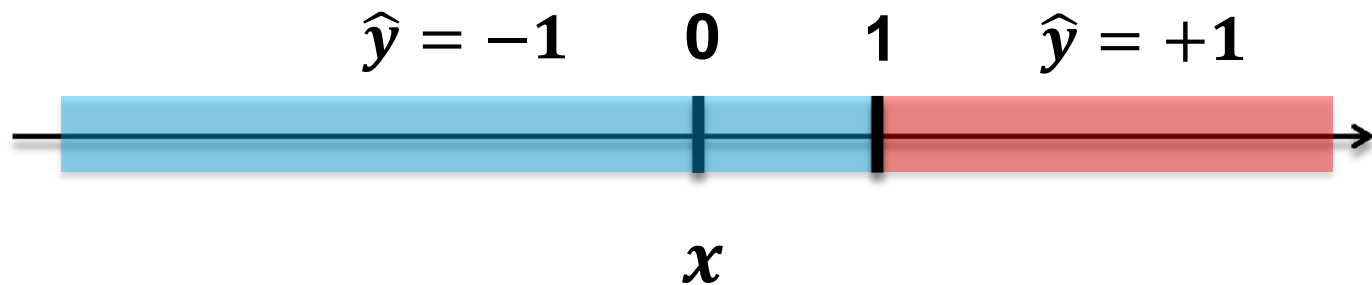
Linear classification

$$\hat{y} = \text{sign}(5 \cdot x - 5)$$



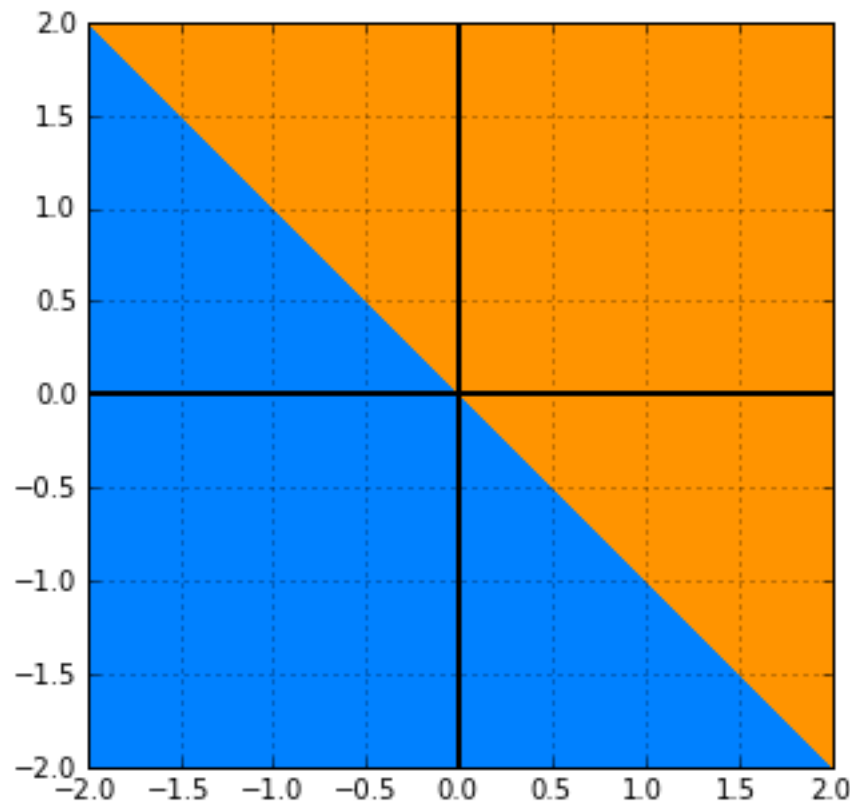
Linear classification

$$\hat{y} = \text{sign}(5 \cdot x - 5)$$
$$\hat{y} = \text{sign}(5 \cdot (x - 1))$$



Linear classification

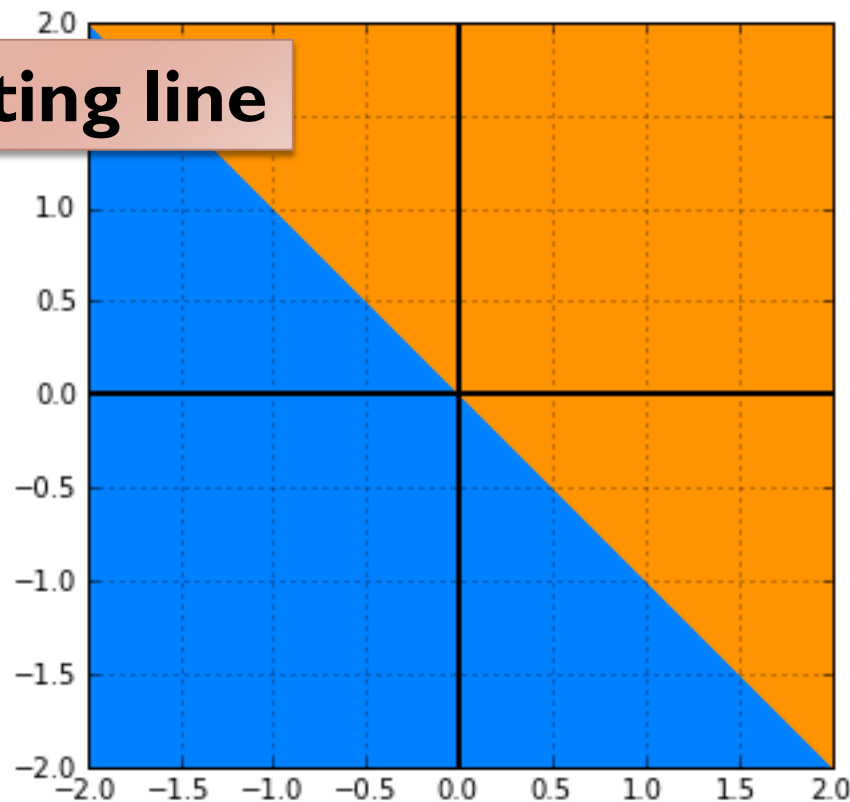
$$\hat{y} = \text{sign}(0.5x_1 + 0.5x_2)$$



Linear classification

$$\hat{y} = \text{sign}(0.5x_1 + 0.5x_2)$$

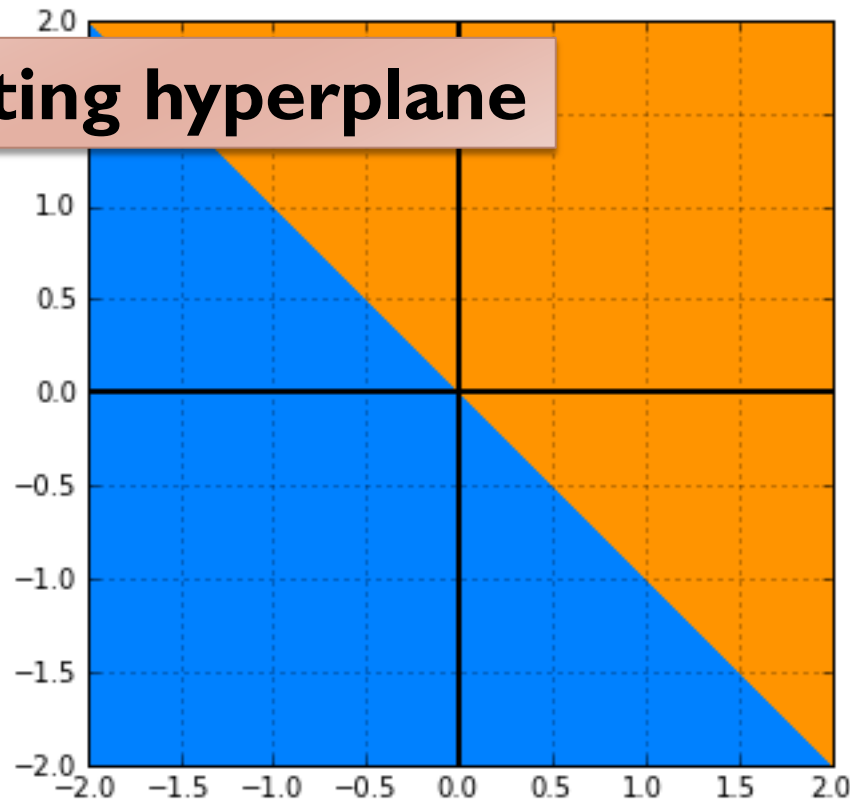
Separating line



Linear classification

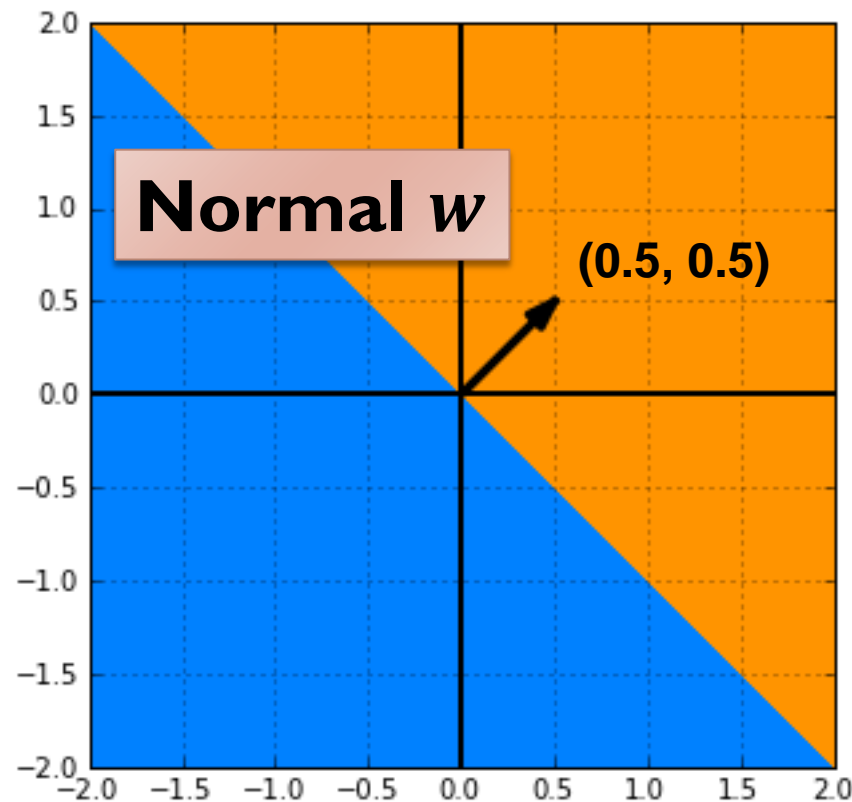
$$\hat{y} = \text{sign}(0.5x_1 + 0.5x_2)$$

Separating hyperplane



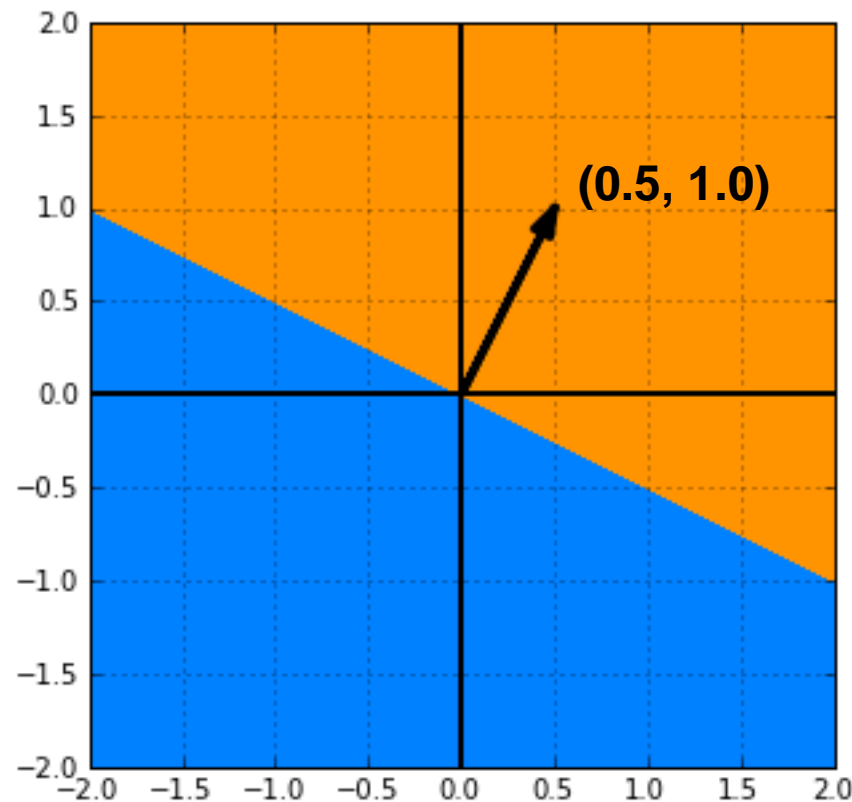
Linear classification

$$\hat{y} = \text{sign}(0.5x_1 + 0.5x_2)$$



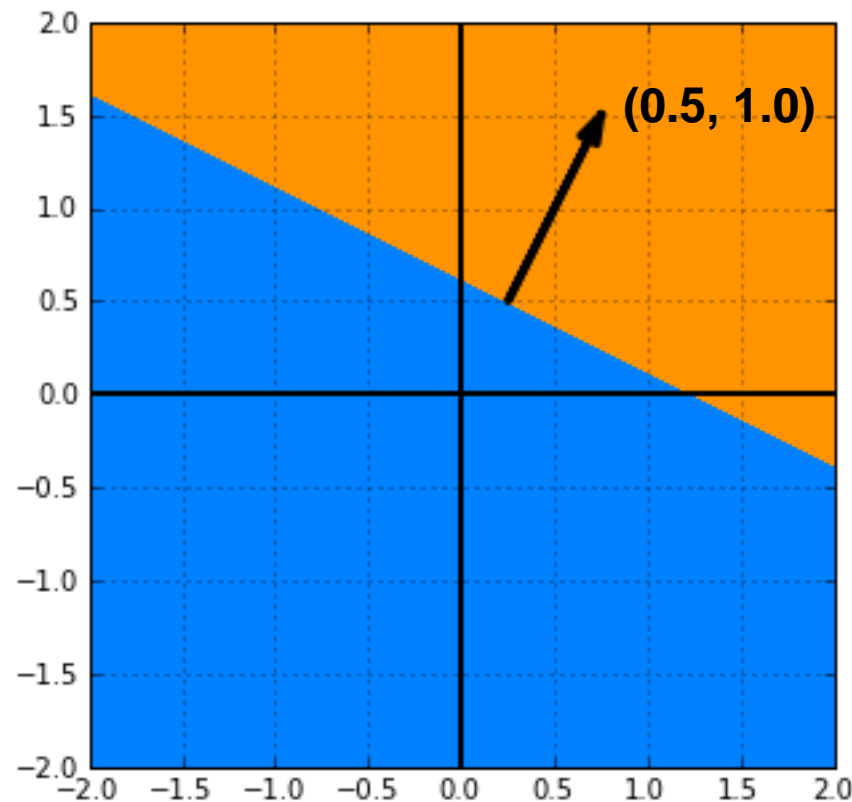
Linear classification

$$\hat{y} = \text{sign}(0.5x_1 + 1.0x_2)$$



Linear classification

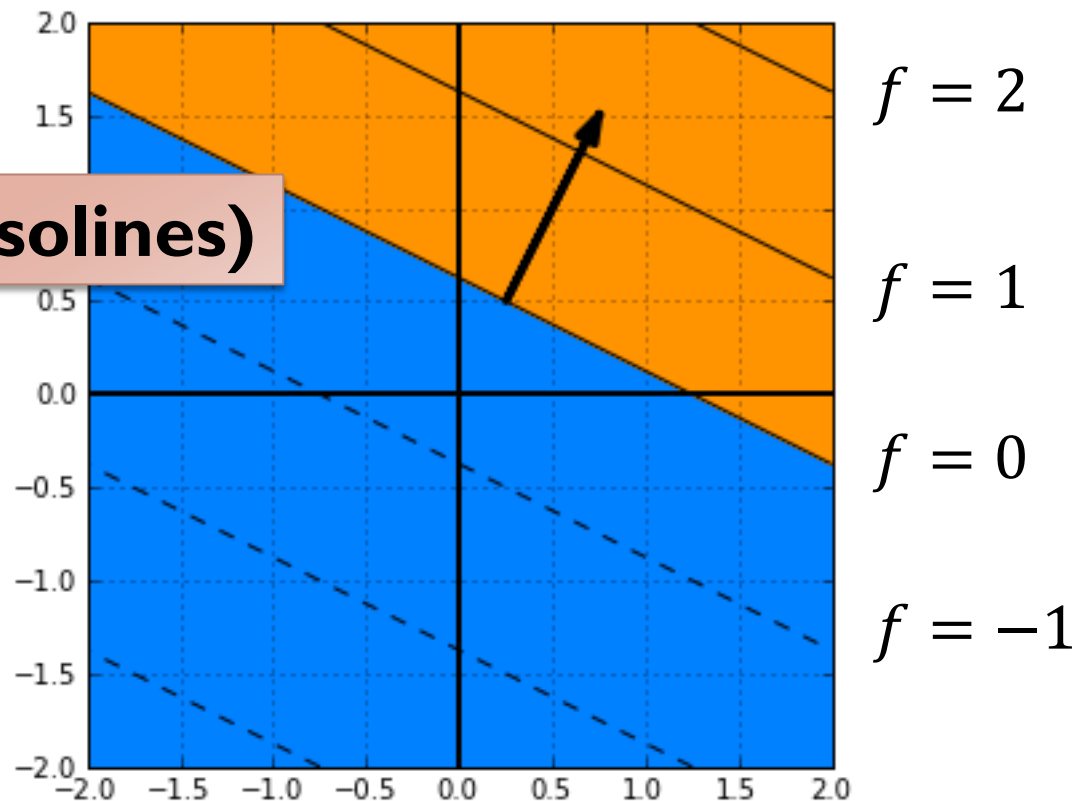
$$\hat{y} = \text{sign}(0.5x_1 + 1.0x_2 - 0.625)$$



Linear classification

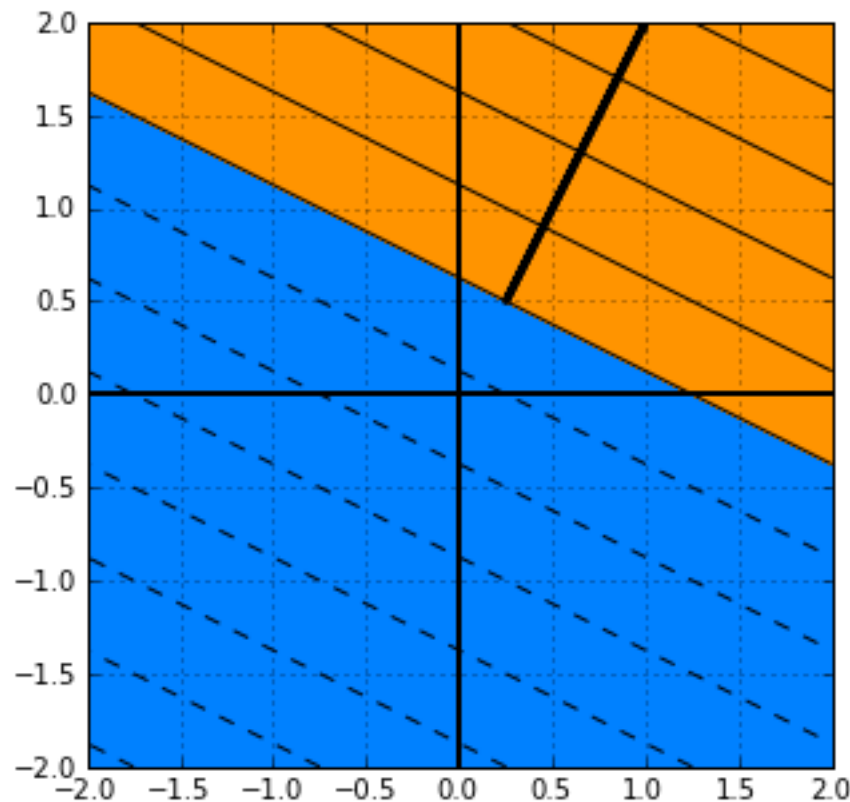
$$\hat{y} = \text{sign}(0.5x_1 + 1.0x_2 - 0.625)$$

Level lines (isolines)



Linear classification

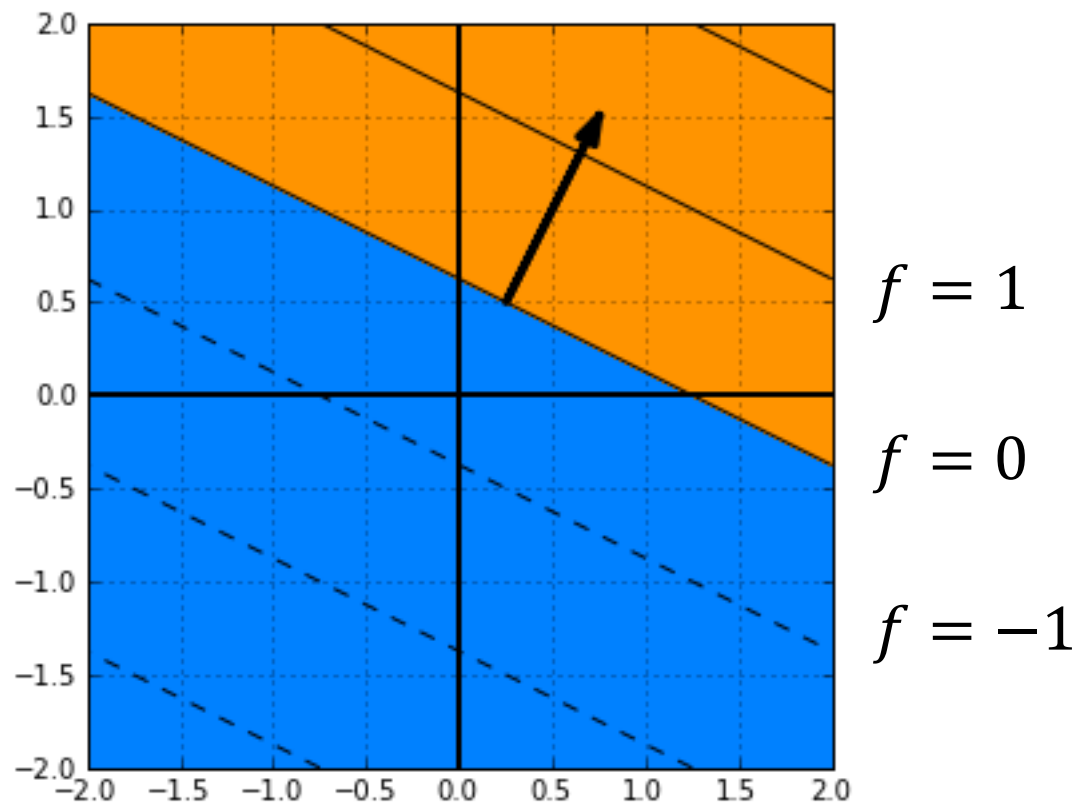
$$\hat{y} = \text{sign}(1.0x_1 + 2.0x_2 - 1.250)$$



$f = 1$
 $f = 0$
 $f = -1$

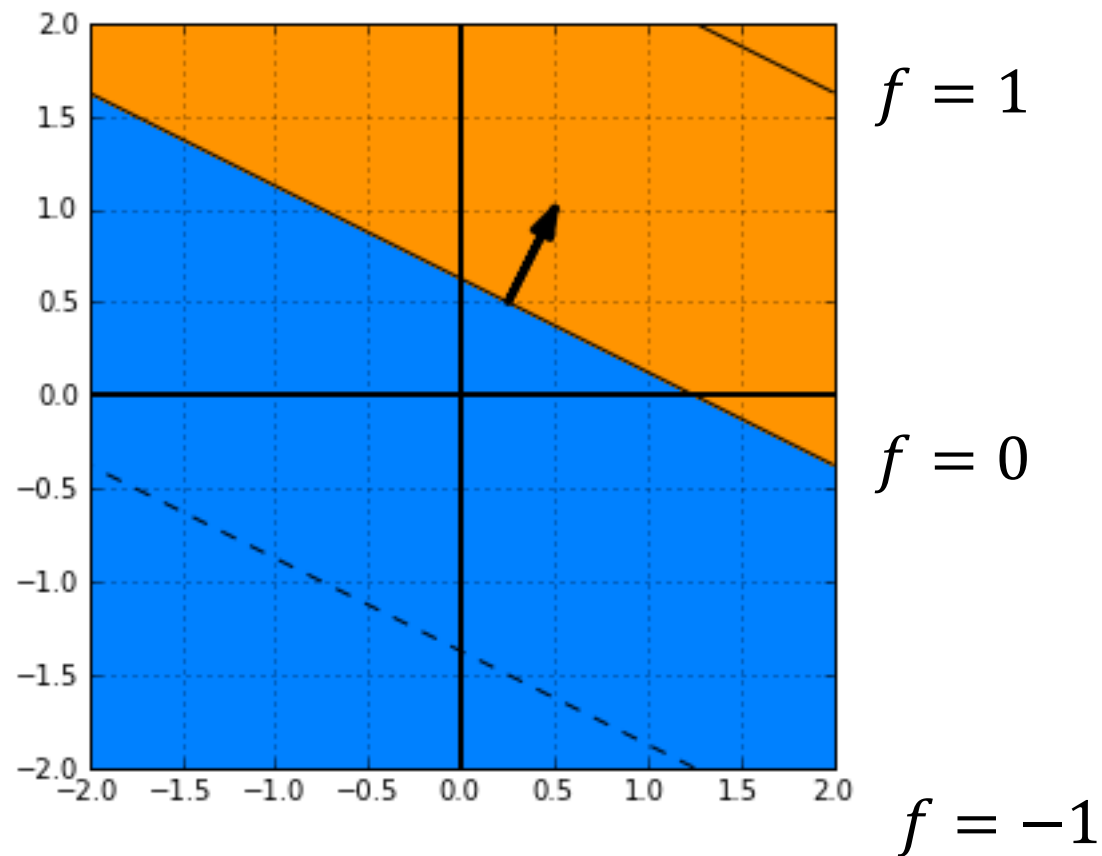
Linear classification

$$\hat{y} = \text{sign}(0.5x_1 + 1.0x_2 - 0.625)$$

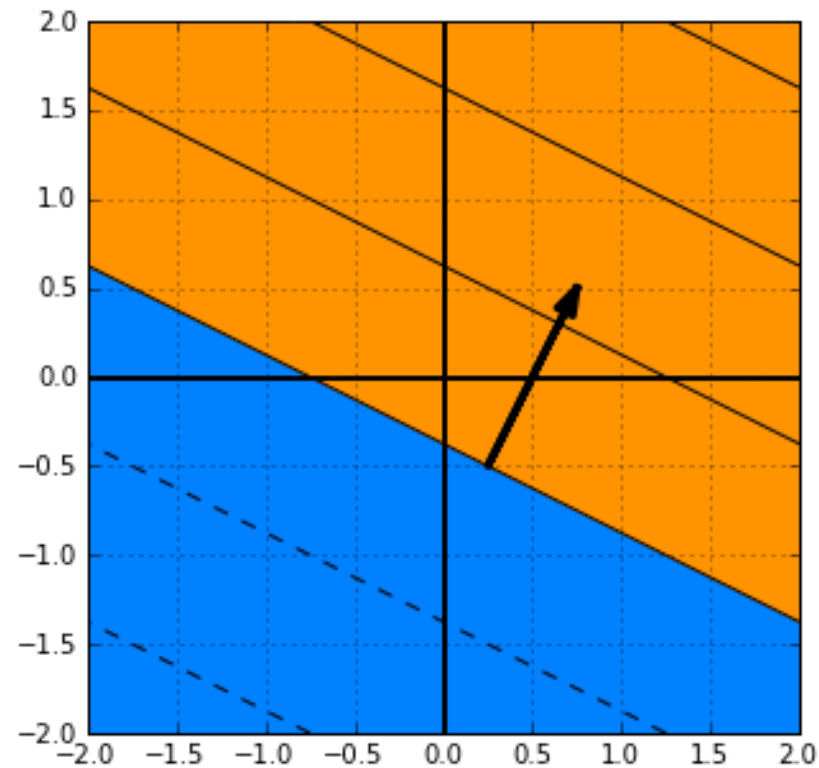


Linear classification

$$\hat{y} = \text{sign}(0.2x_1 + 0.5x_2 - 0.312)$$

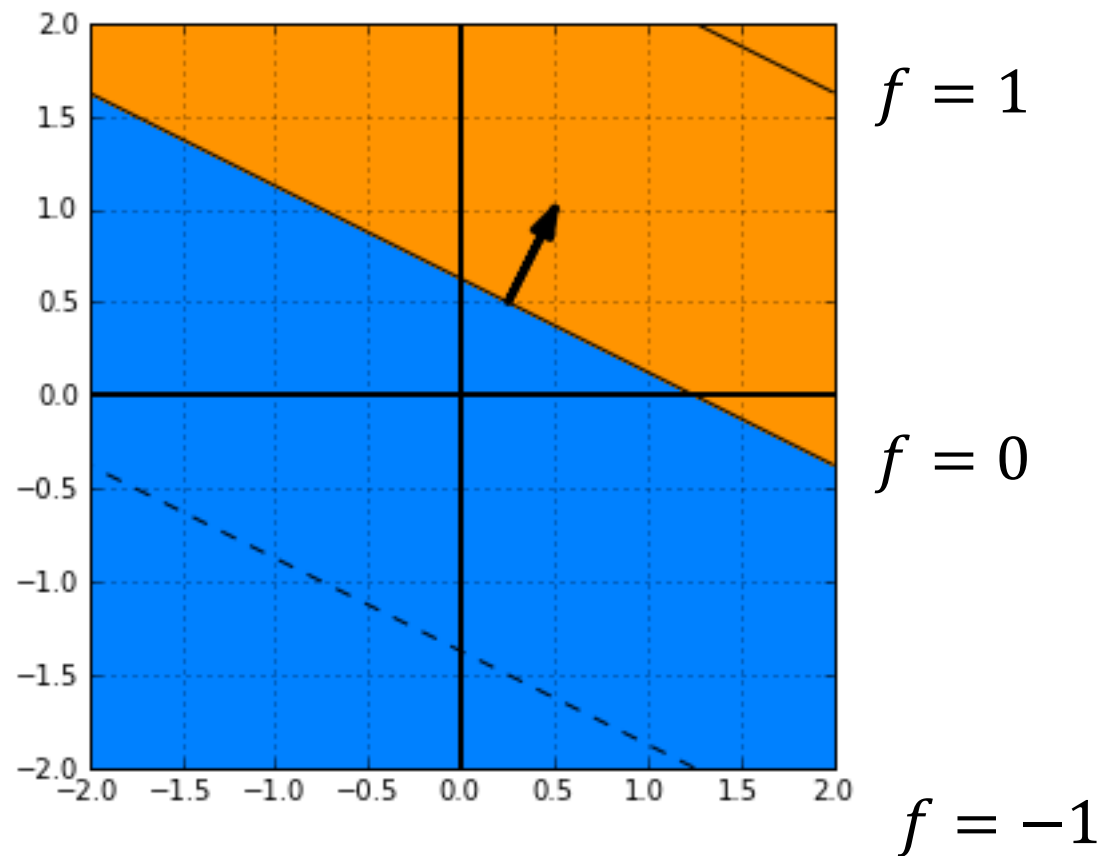


Level Lines



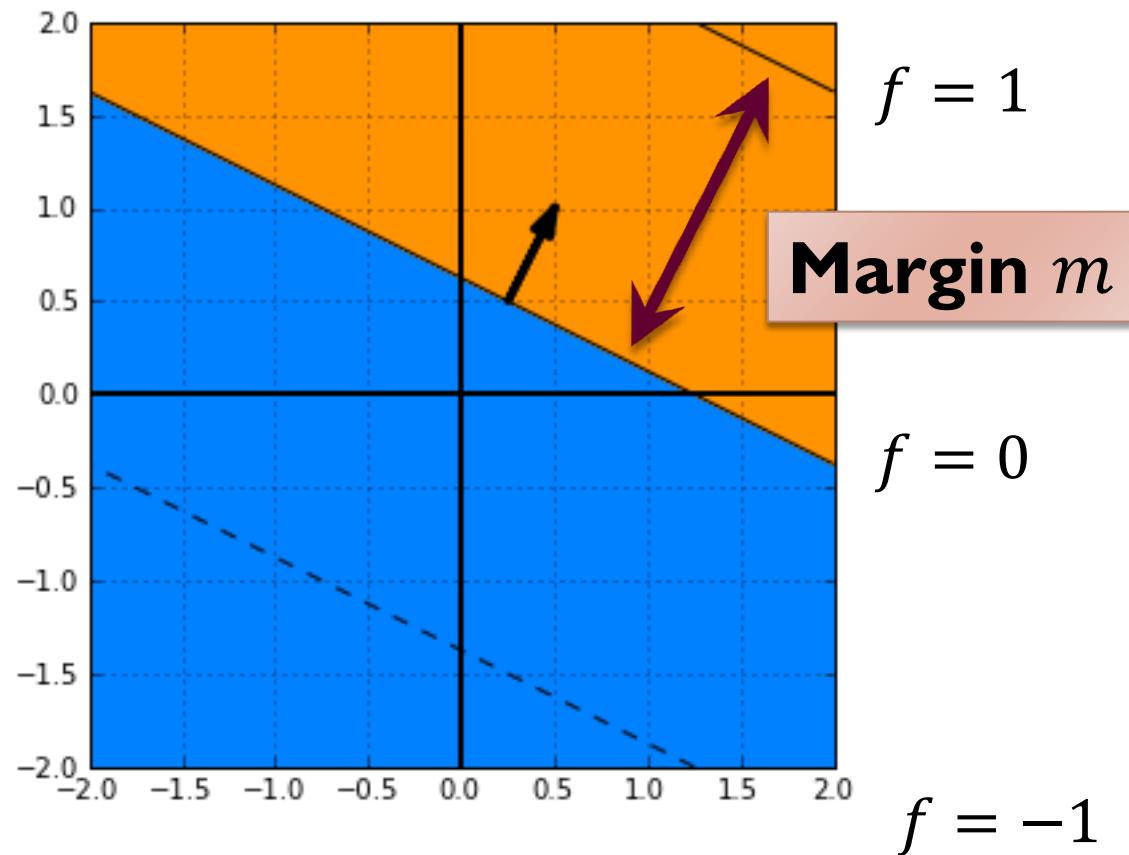
Linear classification

$$\hat{y} = \text{sign}(0.2x_1 + 0.5x_2 - 0.312)$$



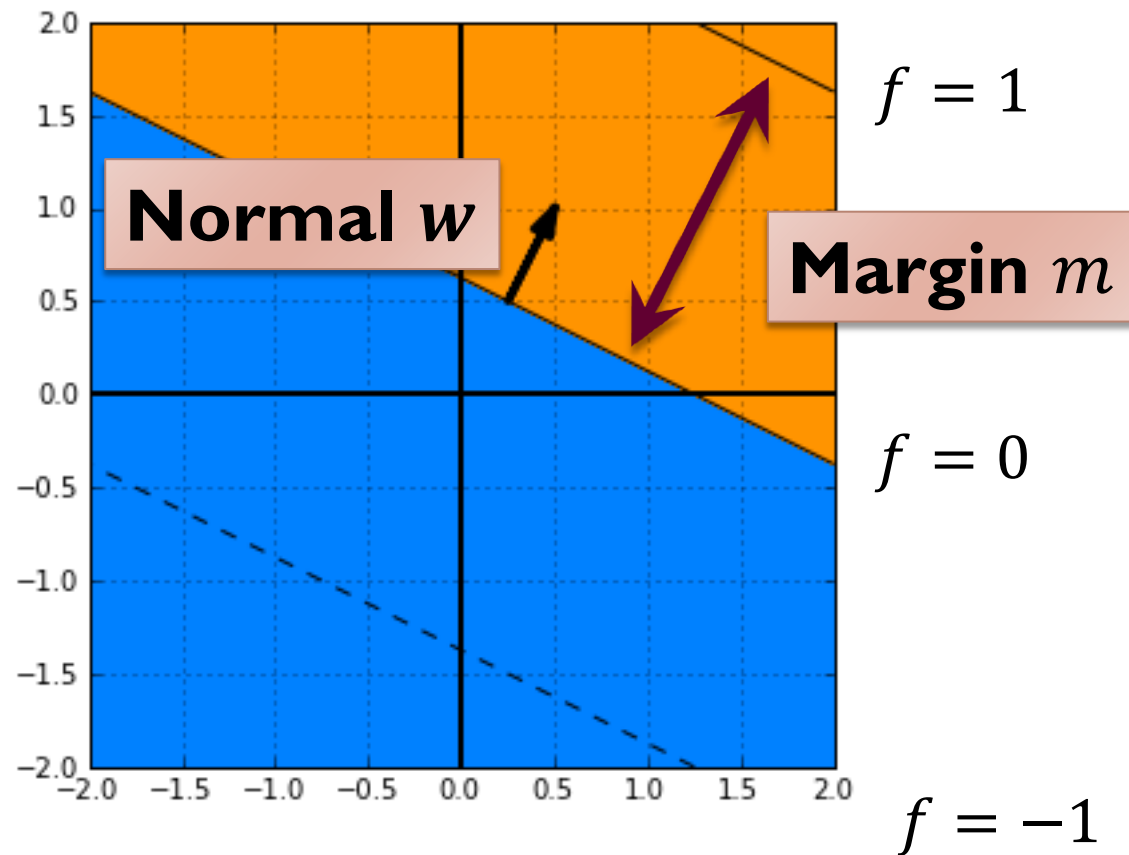
Linear classification

$$\hat{y} = \text{sign}(0.2x_1 + 0.5x_2 - 0.312)$$



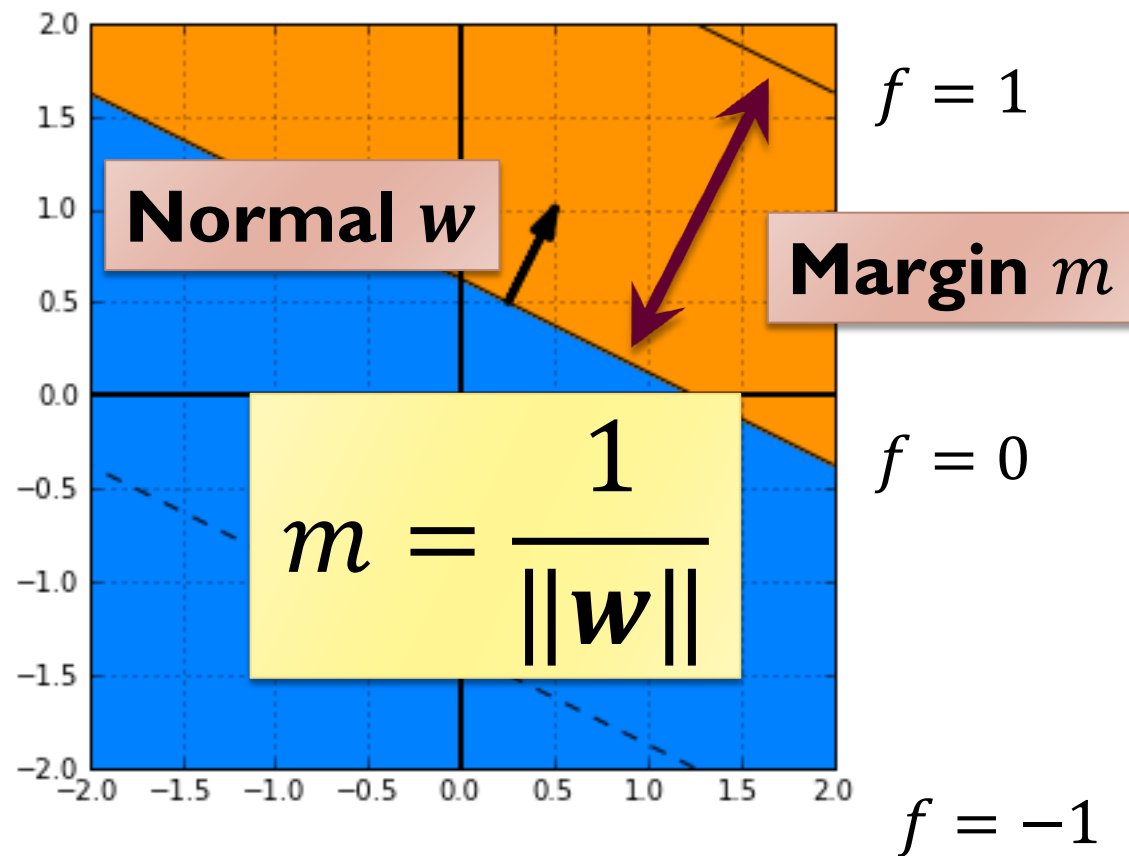
Linear classification

$$\hat{y} = \text{sign}(0.2x_1 + 0.5x_2 - 0.312)$$

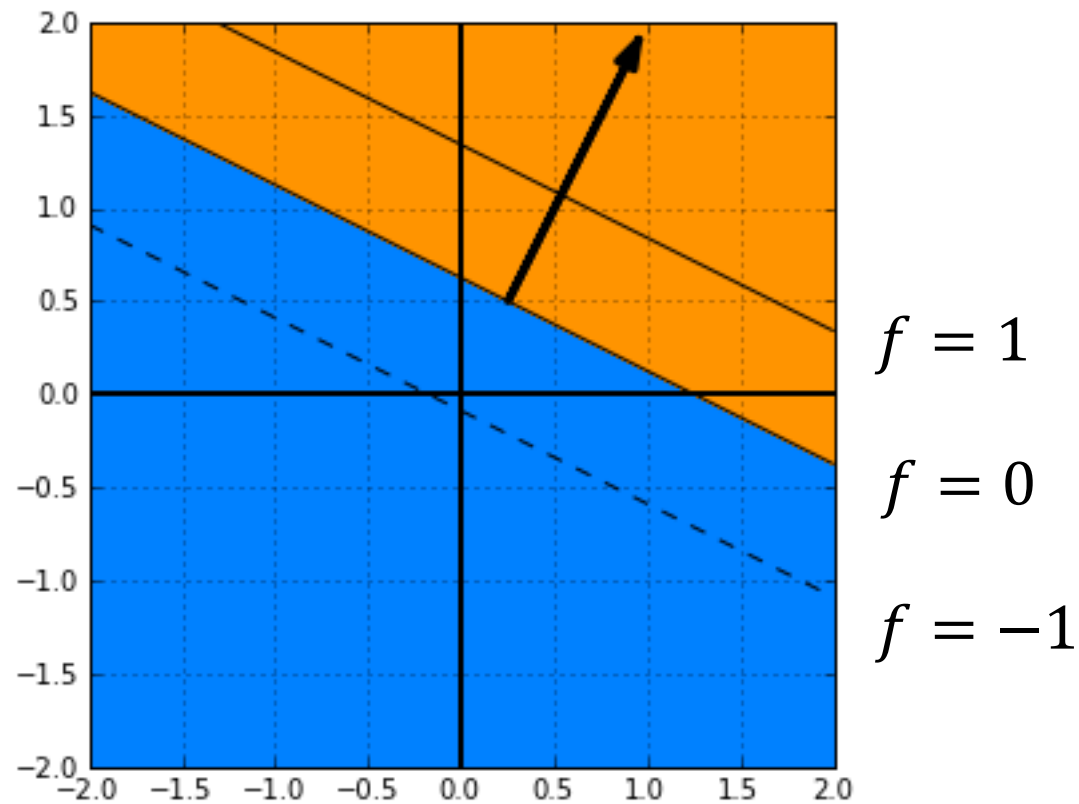


Linear classification

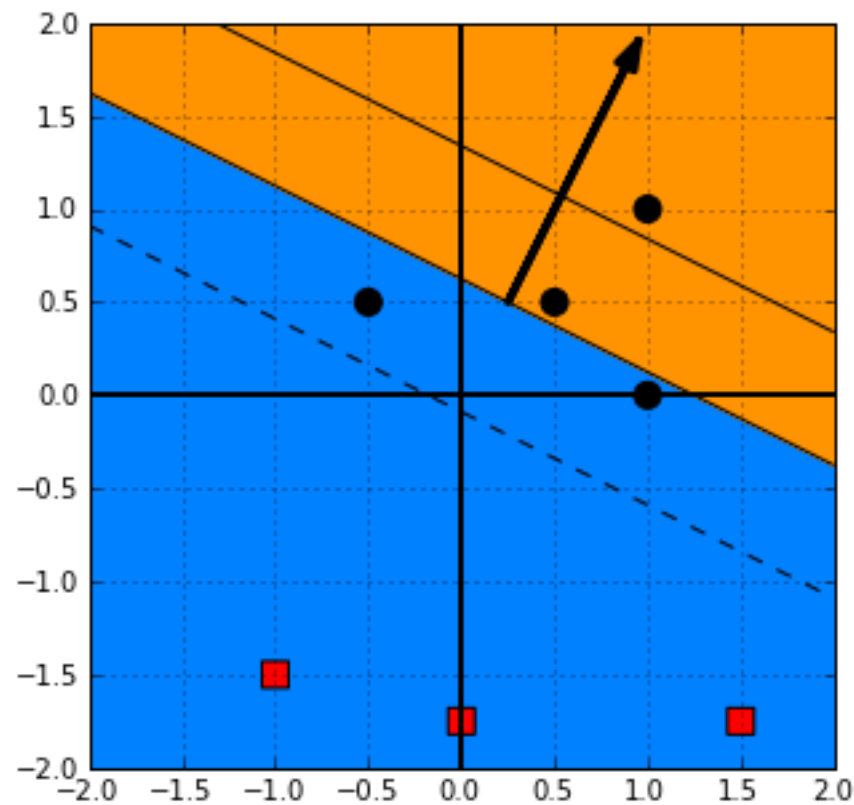
$$\hat{y} = \text{sign}(0.2x_1 + 0.5x_2 - 0.312)$$



Linear classification

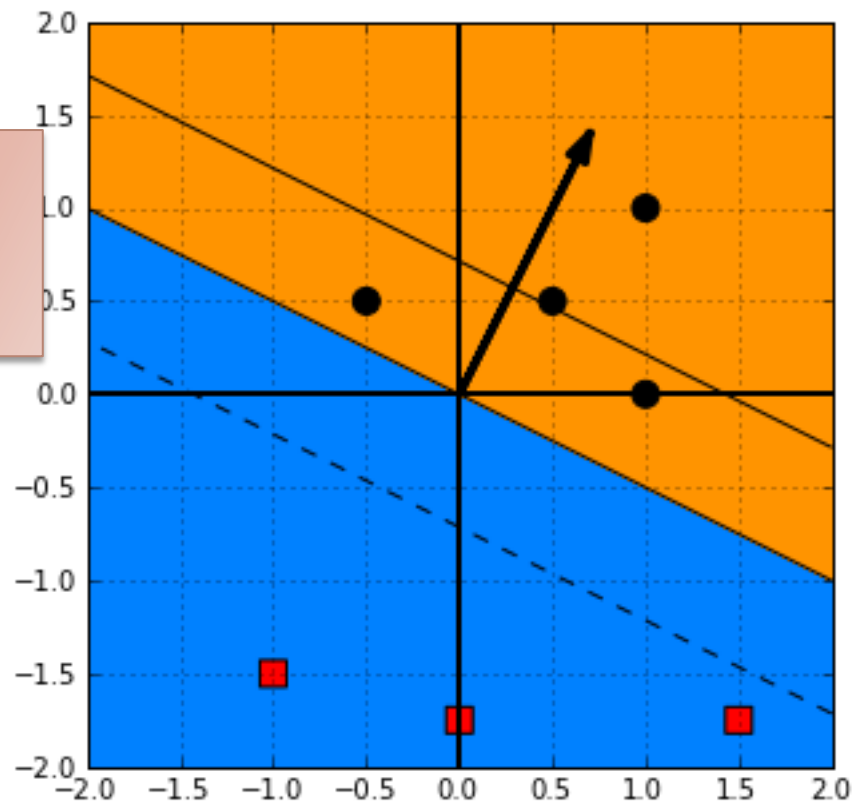


Linear classification



Linear classification

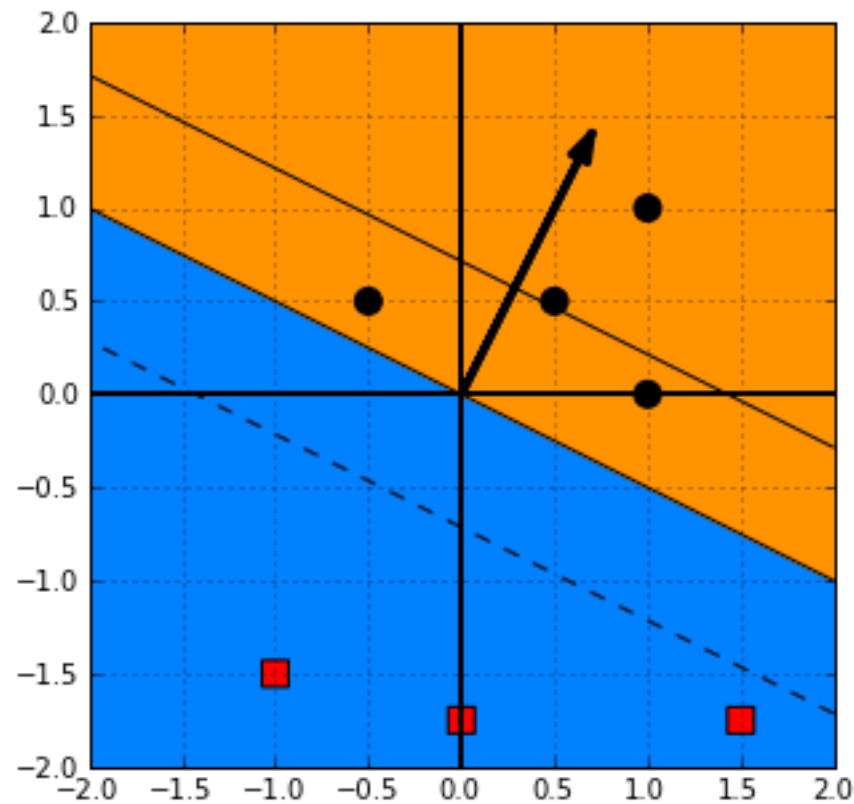
Linearly
separable



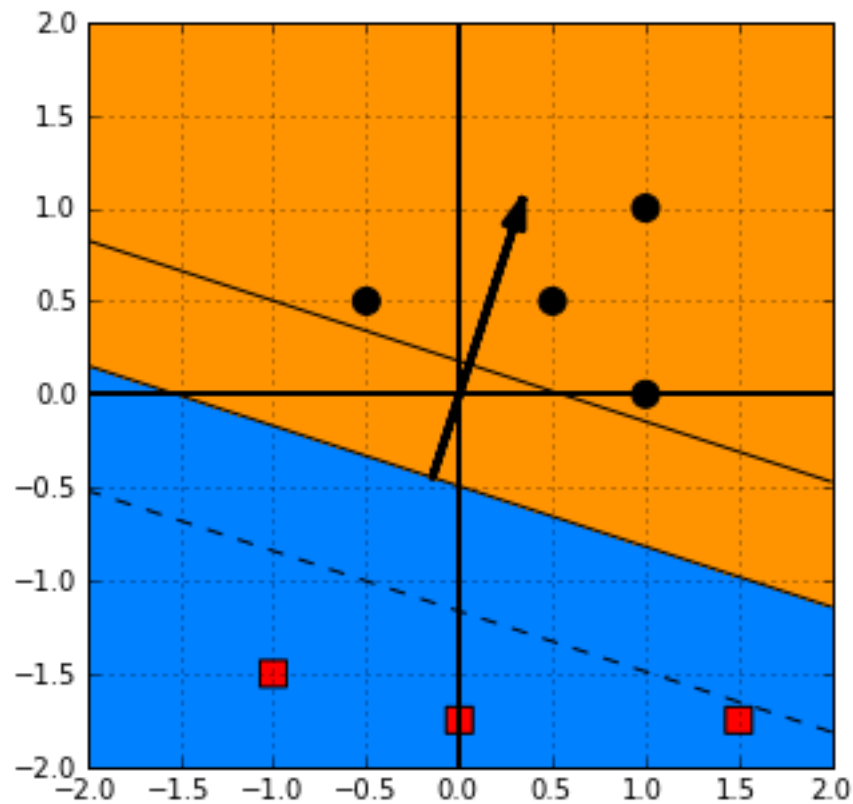
Perceptron

- ▶ Start with arbitrary (\mathbf{w}, w_0)
- ▶ Find a misclassified example (\mathbf{x}_i, y_i)
 - ▶ i.e. $\text{sign}(\mathbf{w}^T \mathbf{x}_i + w_0) \neq y_i$
- ▶ Update:
 - ▶ $\mathbf{w} := \mathbf{w} + y_i \mathbf{x}_i$
 - ▶ $w_0 := w_0 + y_i$

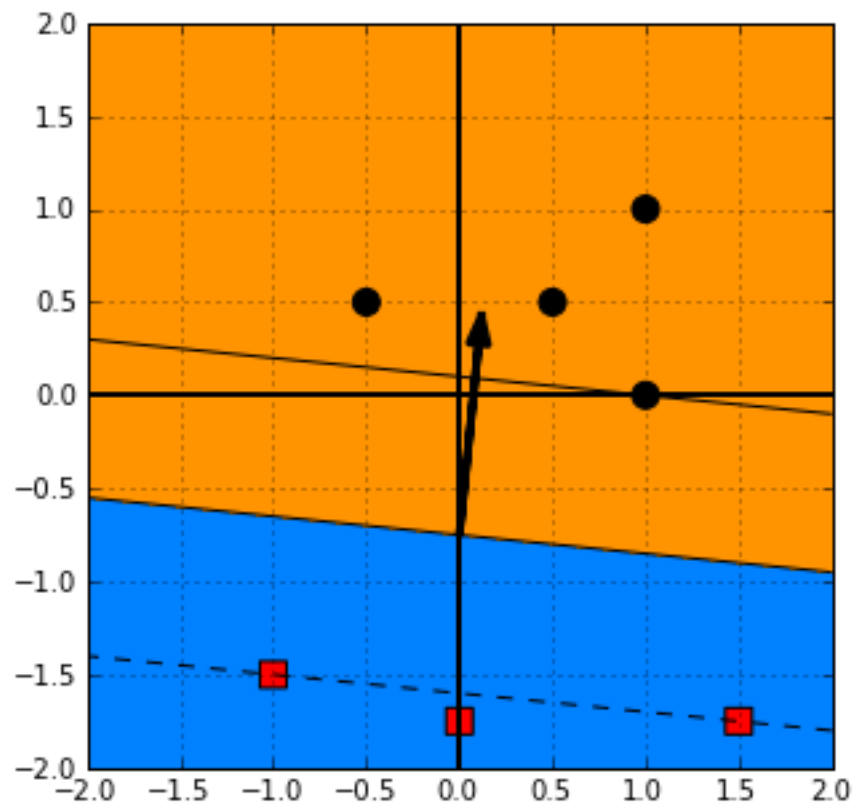
Linear classification



Linear classification



Linear classification

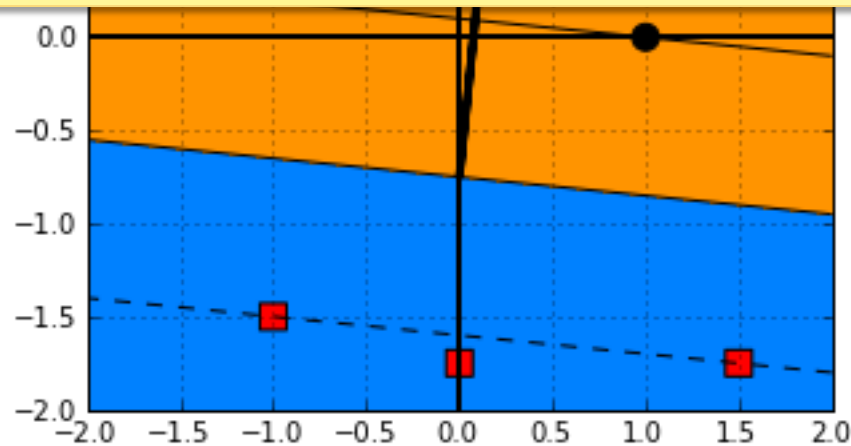


Maximal Margin Classifier

Maximal Margin Classifier

$$\operatorname{argmax}_w m$$

s.t. no points violate the margin

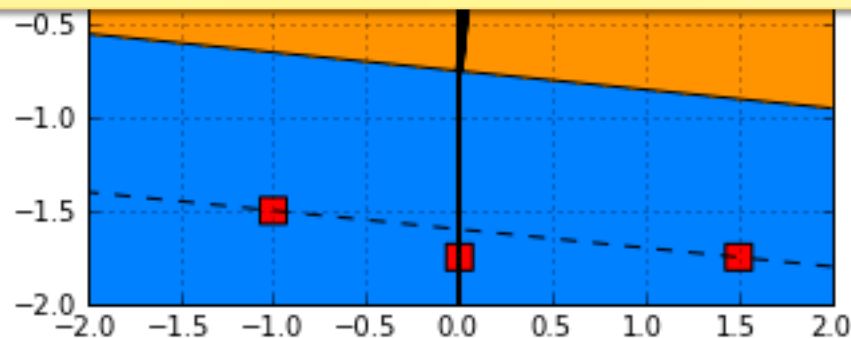


Maximal Margin Classifier

Maximal Margin Classifier

$$\operatorname{argmax}_w \frac{1}{\|w\|}$$

s.t. no points violate the margin

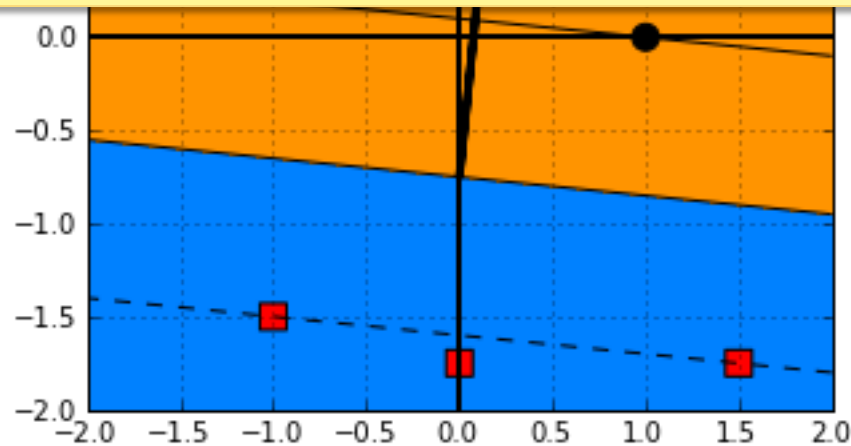


Maximal Margin Classifier

Maximal Margin Classifier

$$\operatorname{argmin}_w \|w\|$$

s.t. no points violate the margin

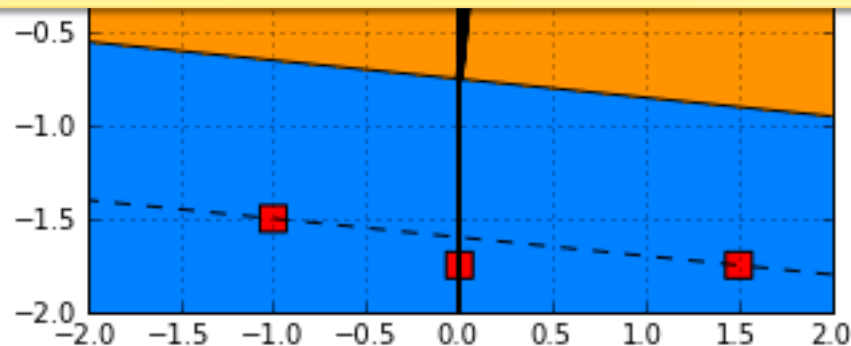


Maximal Margin Classifier

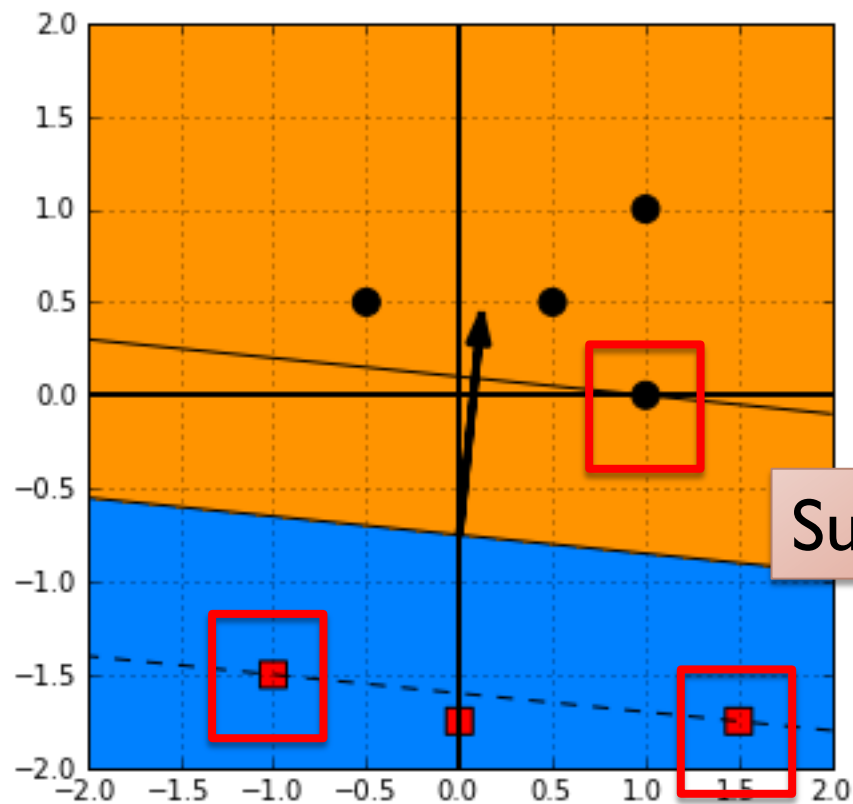
Maximal Margin Classifier

$$\operatorname{argmin}_w \frac{1}{2} \|w\|^2$$

s.t. no points violate the margin

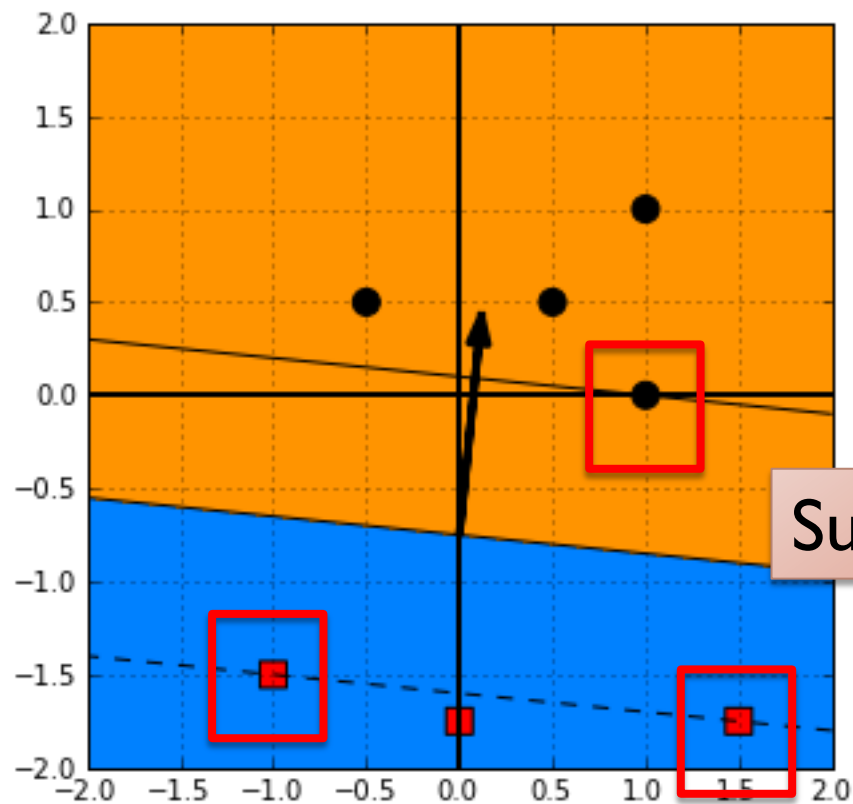


Support Vectors



Support vectors

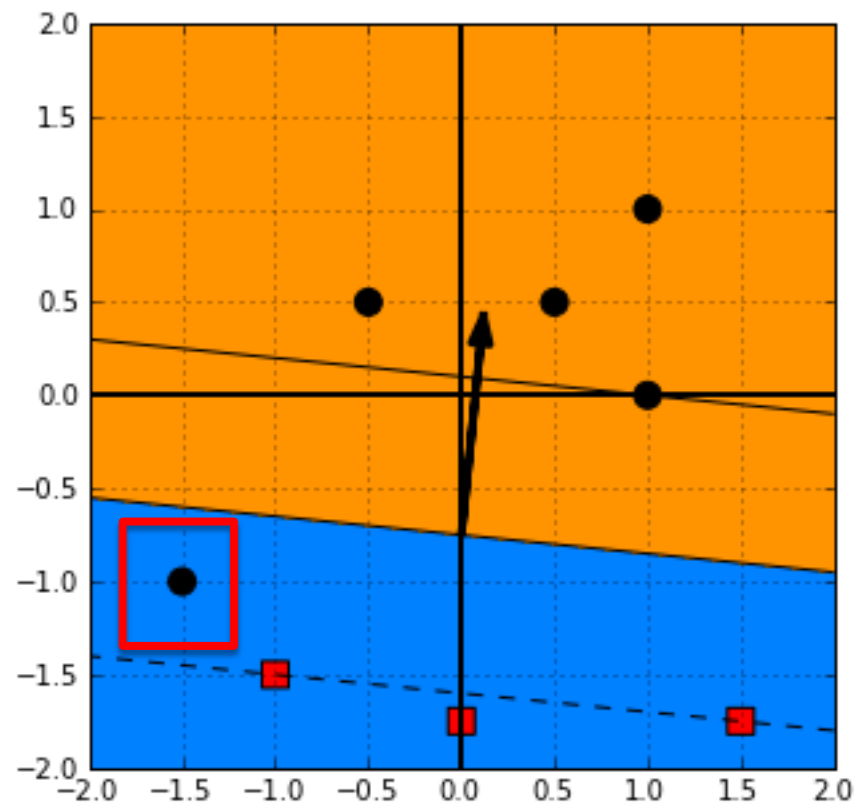
Support Vectors



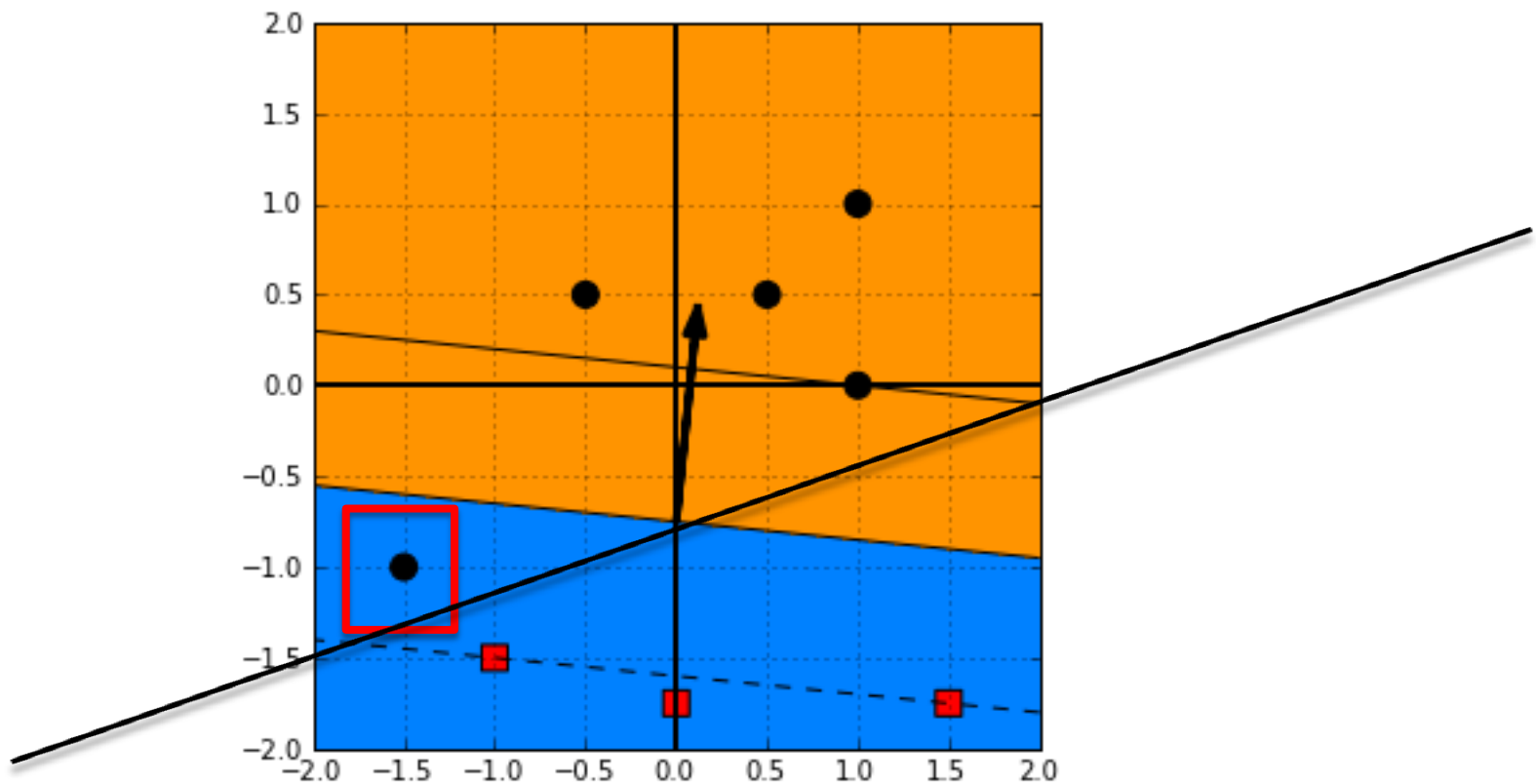
Support vectors

Sparsity!

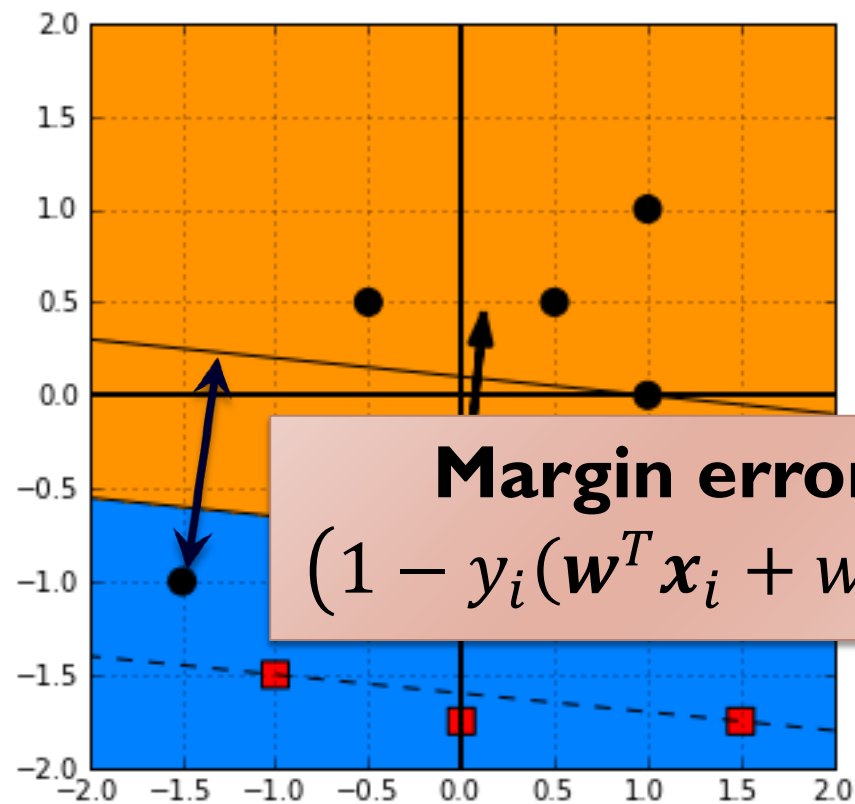
What about outliers?



What about outliers?



Soft-Margin Error



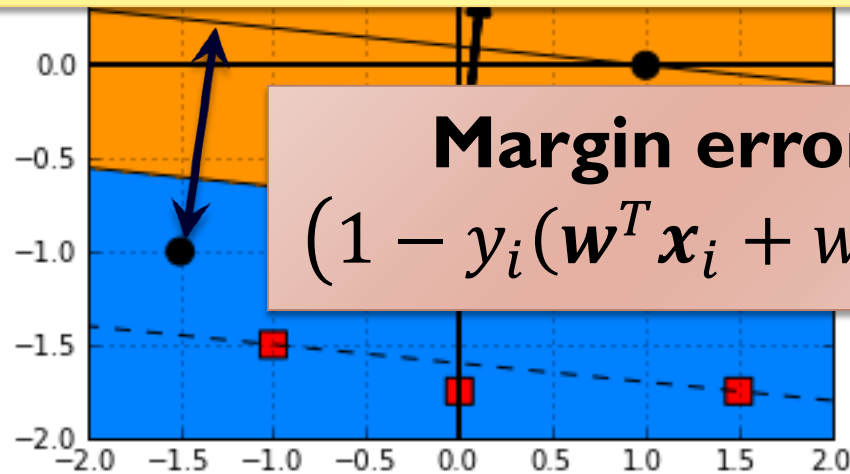
Margin error

$$(1 - y_i(w^T x_i + w_0))_+$$

Soft-Margin Classifier

Soft Margin Classifier

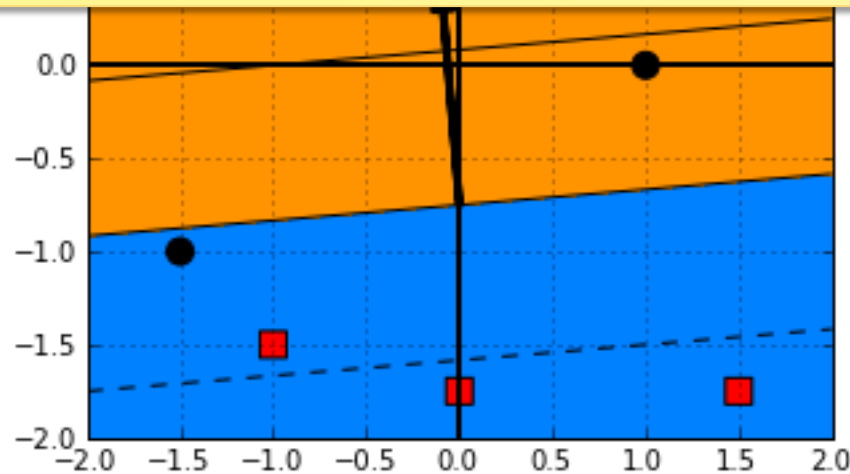
$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))_+$$



Soft-Margin Classifier

Soft Margin Classifier

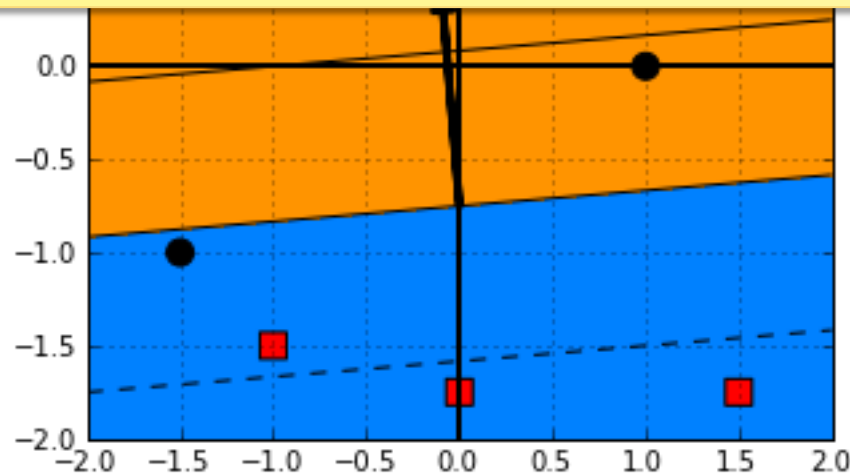
$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))_+$$



Soft-Margin Classifier

Soft Margin Classifier

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))_+$$



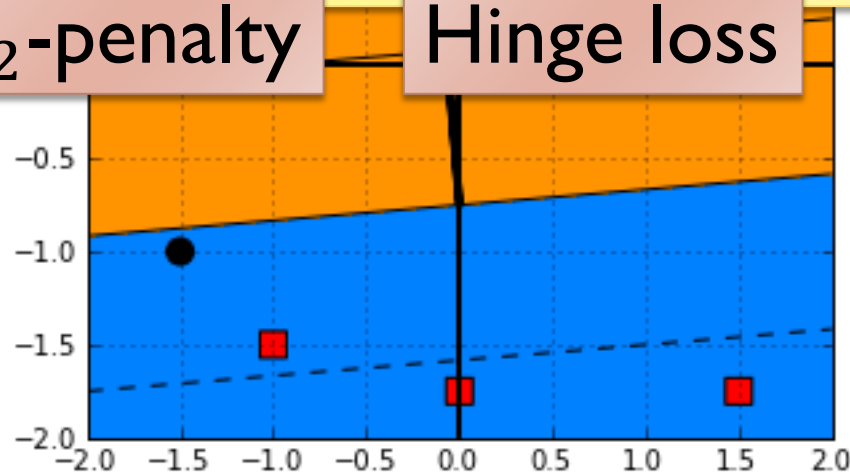
Soft-Margin Classifier

Soft Margin Classifier

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))_+$$

ℓ_2 -penalty

Hinge loss



Support Vector Classifier

```
from sklearn.svm import SVC
```

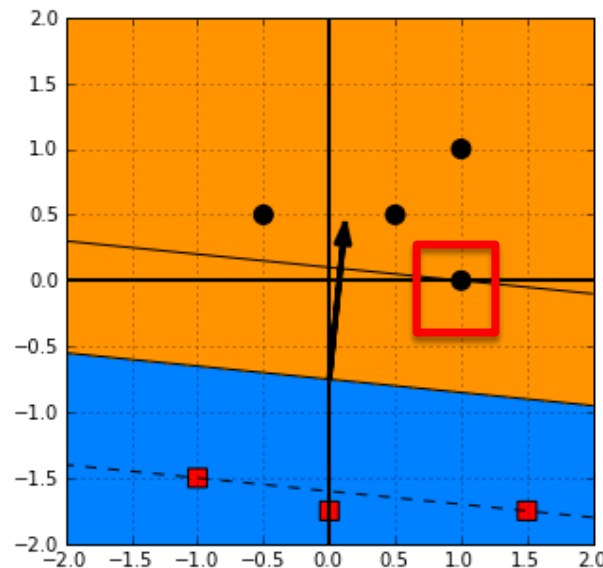
```
m = SVC(C=1, kernel='linear')
```

```
m.fit(X, y)
```

```
m.predict(X_new)
```

Quiz

- ▶ The smaller is the classifier's _____, the larger is its margin.
- ▶ If the point lies exactly on the margin (on the correct side), it is called _____.



Quiz

- ▶ If you remove all non-support vectors from the dataset and re-train the model, the resulting hyperplane will be _____.

Quiz

- ▶ The margin error of a point exactly on the margin is _____.
- ▶ The margin error of a point exactly on the separating hyperplane is _____.

Margin error

$$(1 - y_i(\mathbf{w}^T \mathbf{x}_i + w_0))_+$$

Next

Reasoning by
analogy explains
it all!

It is plain old
statistics and
probability
theory!

Just a bunch of
useful methods
and tools.

It is a branch of
cognitive science!



Linear classification & SVMs

Kernel Methods

an expanded
version of the
specific method

The Kernel-based Approach

Nonlinear feature mapping

The Kernel trick

Dual representation



-
- ▶ Linear regression (OLS, Ridge):

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- ▶ Linear classification (SVM):

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0)$$

- ▶ Linear regression (OLS, Ridge, LASSO, SVR,...):

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- ▶ Linear classification (SVM, Perceptron, Fisher's discriminant, Logistic regression, NB, ...):

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0)$$

- ▶ Linear regression (OLS, Ridge, LASSO, SVR,...):

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- ▶ Linear classification (SVM, Perceptron, Fisher's discriminant, Logistic regression, NB, ...):

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0)$$

- ▶ PCA, LDA, ICA, CCA...:

$$\mathbf{x}_T = \mathbf{A}\mathbf{x}$$

- ▶ K-means:

$$\mathbf{c}_i = \frac{1}{m} \mathbf{X}_i \mathbf{1}$$

▶ Linear regression (OLS, Ridge, LASSO, SVR,...):

▶ Line
disc

▶ PCA

▶ K-m

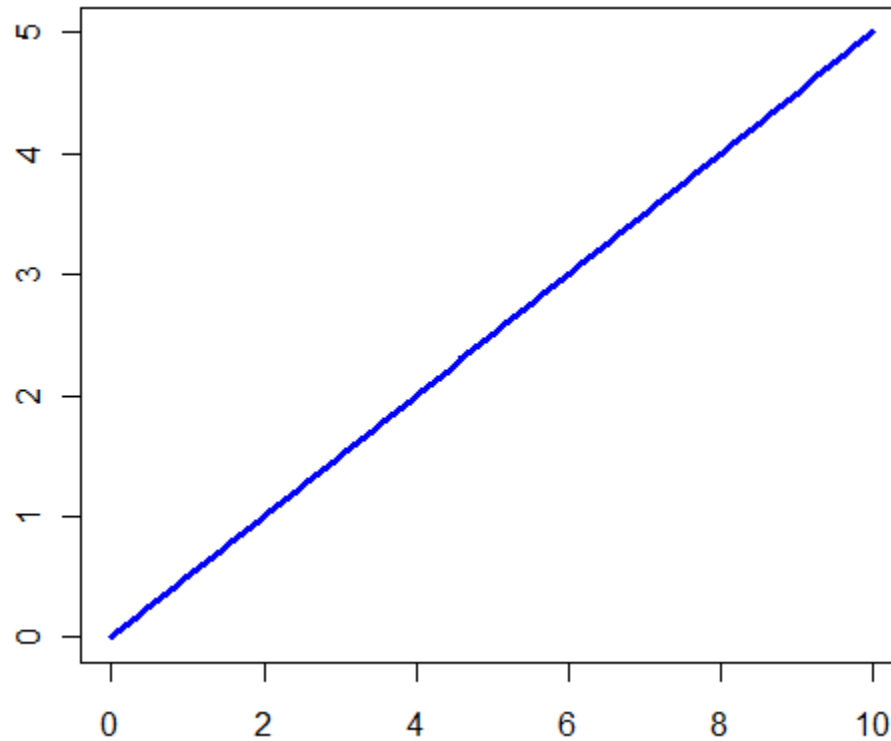


her's

$$c_i = \frac{1}{m} \mathbf{X}_i \mathbf{1}$$

Recall polynomial regression

$$f(x) = wx$$



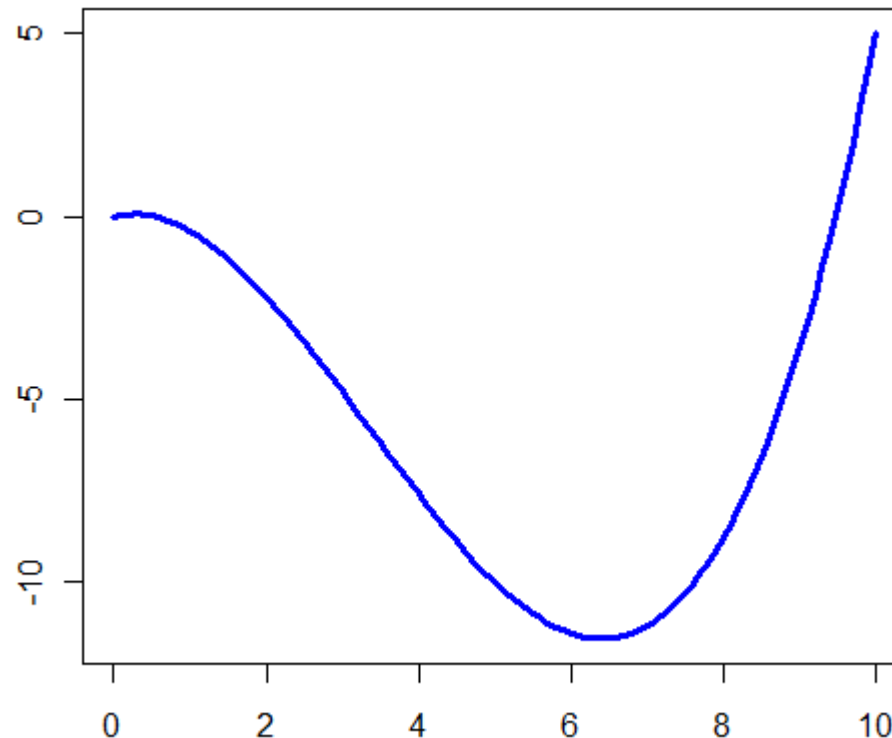
Recall polynomial regression

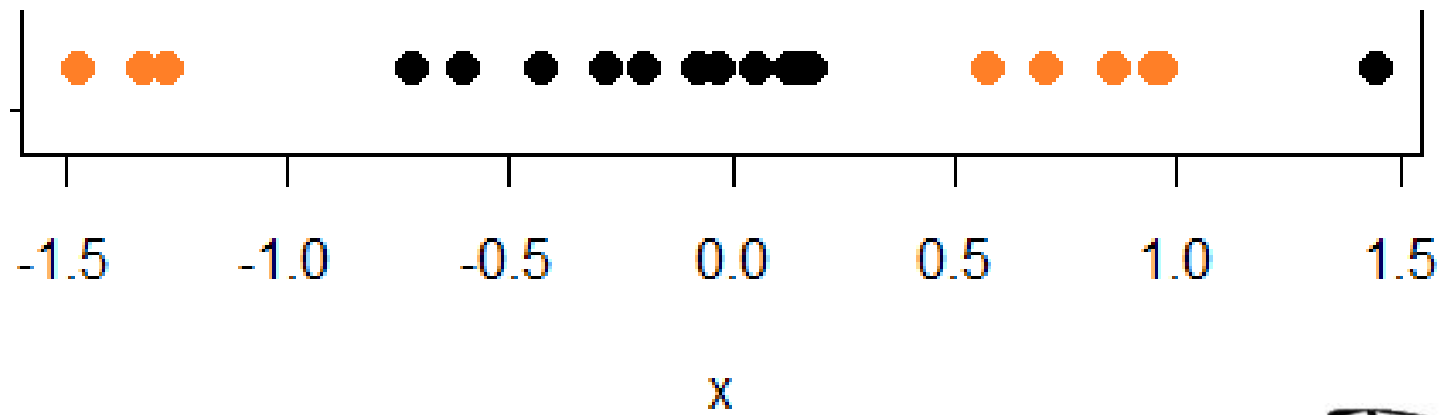
$$x \rightarrow x' := \phi(x) := (x, x^2, x^3)$$

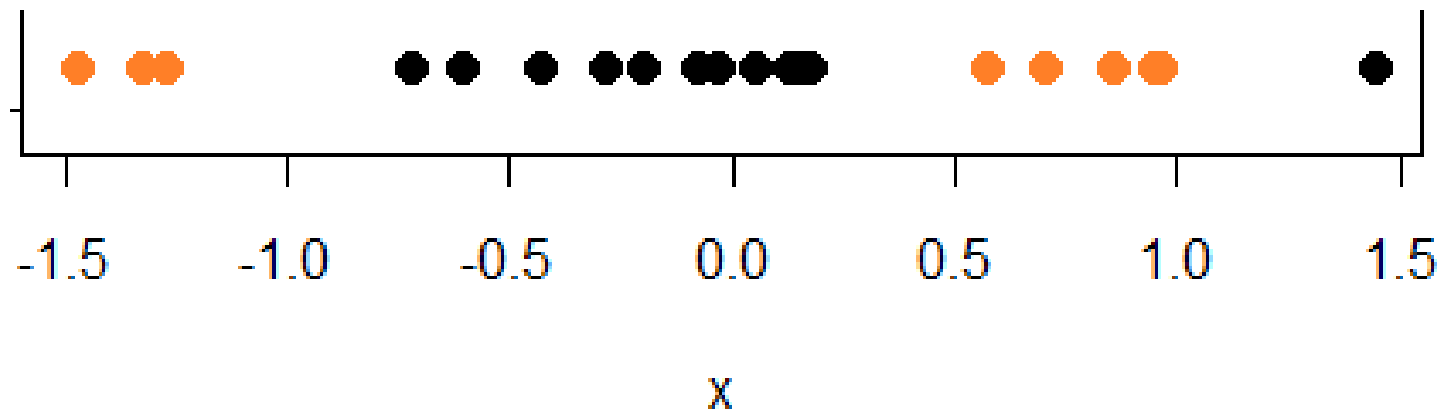


Nonlinear feature space

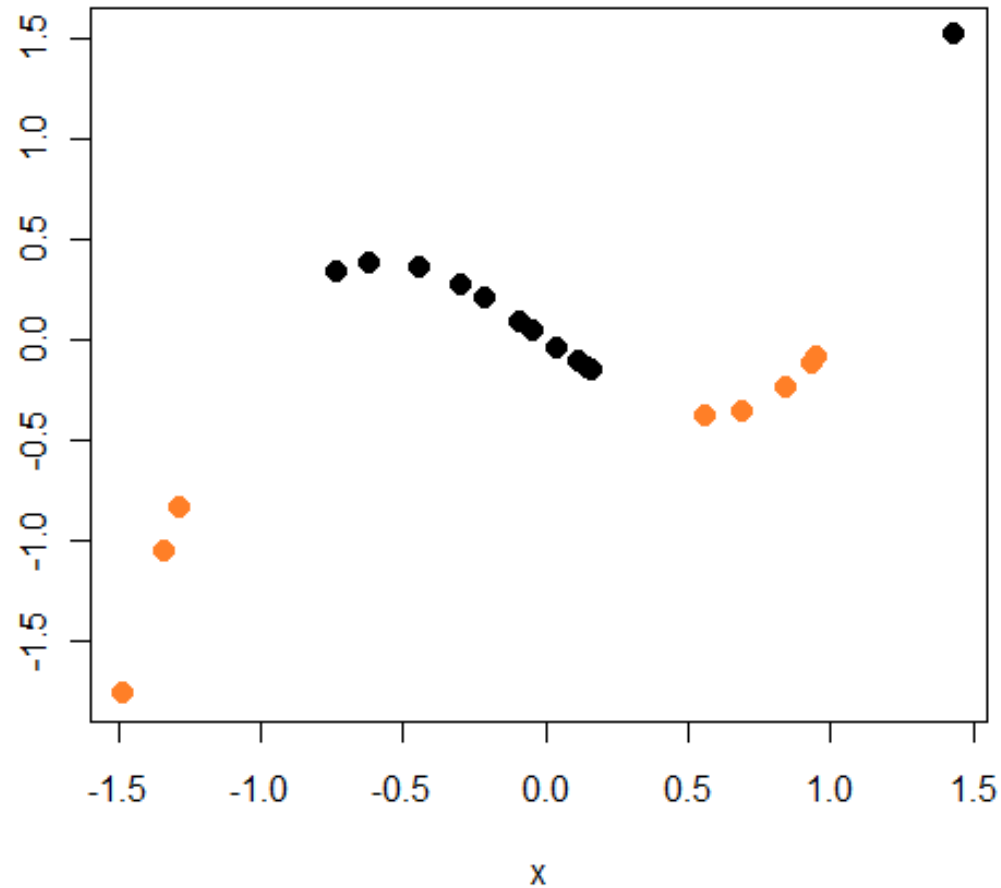
$$x \rightarrow x' := \phi(x) := (x, x^2, x^3)$$
$$f(x') = w_1x + w_2x^2 + w_3x^3$$



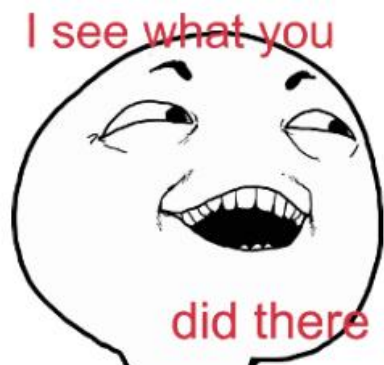
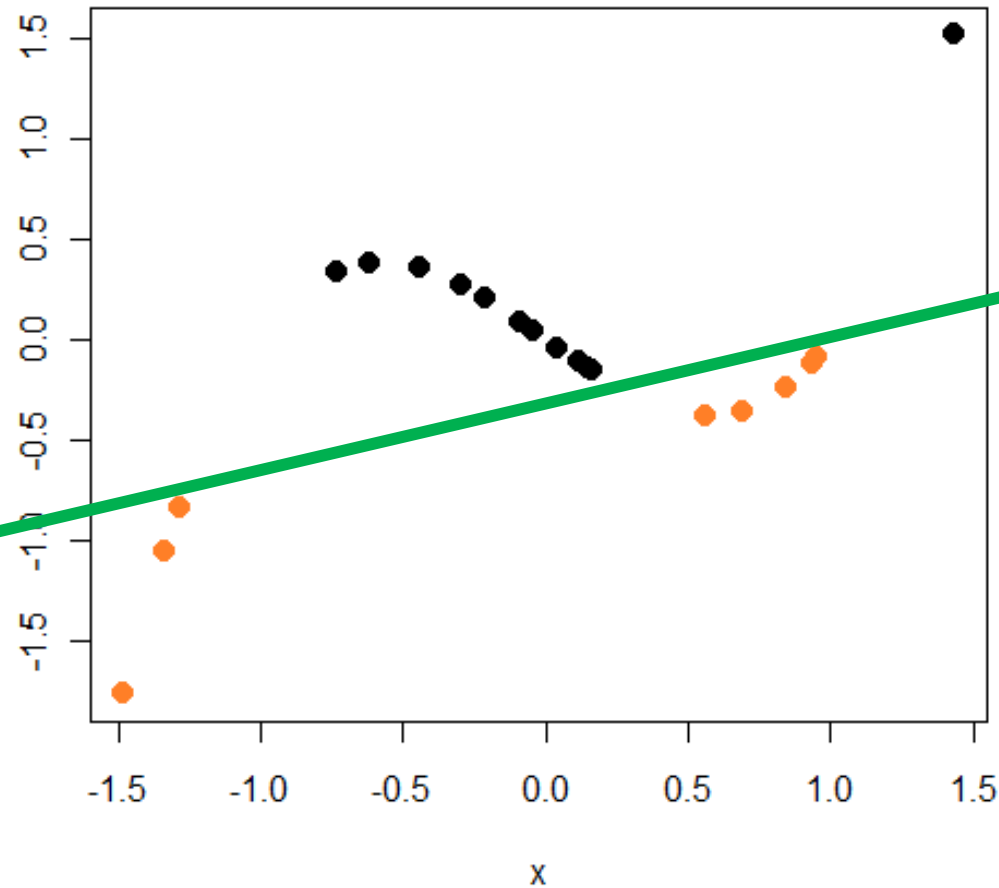




$$x \rightarrow \phi(x) = (x, x^3 - x)$$



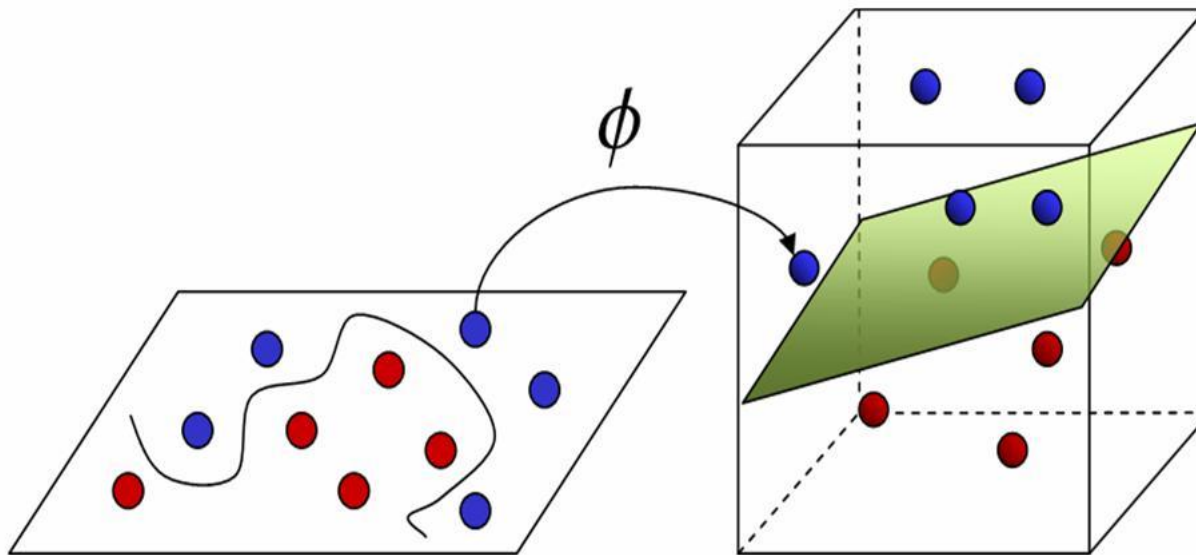
$$x \rightarrow \phi(x) = (x, x^3 - x)$$



$$x \rightarrow \phi(x) = (x, x^3 - x)$$

Nonlinear feature space

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$



Support for arbitrary data types

$\phi(\text{text}) = \text{word counts}$

$\phi(\text{graph}) = \text{node degrees}$

$\phi(\text{tree}) = \text{path lengths}$

...

What if the dimensionality is high?

$$(x_1, x_2, \dots, x_m) \rightarrow (x_1 x_1, x_1 x_2, \dots, x_m x_m)$$

What if the dimensionality is high?

$$(x_1, x_2, \dots, x_m) \rightarrow (x_1x_1, x_1x_2, \dots, x_mx_m)$$

$O(m^2)$ elements

For all k-wise products: $O(m^k)$



The Kernel-based Approach

Nonlinear feature mapping

The Kernel trick

Dual representation



The Kernel Trick

- ▶ Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$
- ▶ Consider

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \sum_{ij} \phi(\mathbf{x})_{ij} \phi(\mathbf{y})_{ij}$$

The Kernel Trick

- ▶ Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$
- ▶ Consider

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle &= \sum_{ij} \phi(\mathbf{x})_{ij} \phi(\mathbf{y})_{ij} \\ &= \sum_{ij} x_i x_j y_i y_j\end{aligned}$$

The Kernel Trick

- ▶ Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$
- ▶ Consider

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle &= \sum_{ij} \phi(\mathbf{x})_{ij} \phi(\mathbf{y})_{ij} \\ &= \sum_{ij} x_i x_j y_i y_j = \sum_{ij} x_i y_i x_j y_j\end{aligned}$$

The Kernel Trick

- ▶ Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$
- ▶ Consider

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle &= \sum_{ij} \phi(\mathbf{x})_{ij} \phi(\mathbf{y})_{ij} \\ &= \sum_{ij} x_i x_j y_i y_j = \sum_{ij} x_i y_i x_j y_j \\ &= \sum_i x_i y_i \sum_j x_j y_j\end{aligned}$$

The Kernel Trick

- ▶ Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$
- ▶ Consider

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle &= \sum_{ij} \phi(\mathbf{x})_{ij} \phi(\mathbf{y})_{ij} \\ &= \sum_{ij} x_i x_j y_i y_j = \sum_{ij} x_i y_i x_j y_j \\ &= \sum_i x_i y_i \sum_j x_j y_j = \left(\sum_i x_i y_i \right)^2\end{aligned}$$

The Kernel Trick

► Let $\phi(\mathbf{x}) = (x_1x_1, x_1x_2, \dots, x_mx_m)$



$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \langle \mathbf{x}, \mathbf{y} \rangle^2$$

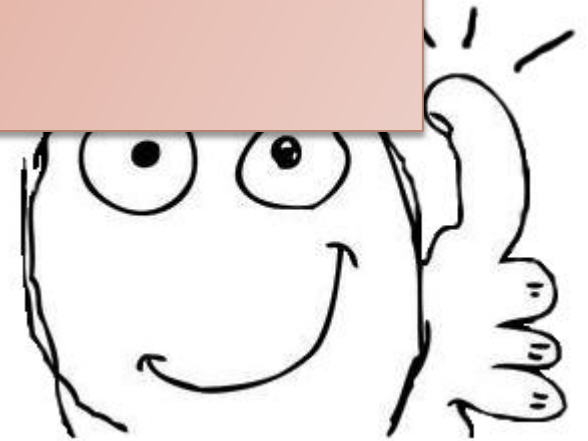


The Kernel Trick

► Let

Polynomial kernel

$$K(x, y) = (\langle x, y \rangle + R)^d$$



The Kernel Trick

What about:

$$K(x, y) = \langle x, y \rangle + 0.5 \langle x, y \rangle^2?$$

The Kernel Trick

What about:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{x}, \mathbf{y} \rangle + 0.5 \langle \mathbf{x}, \mathbf{y} \rangle^2 \\ &= \sum_i x_i y_i + 0.5 \sum_{ij} \phi_{ij}(\mathbf{x}) \phi_{ij}(\mathbf{y}) \end{aligned}$$

The Kernel Trick

What about:

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle + 0.5 \langle \mathbf{x}, \mathbf{y} \rangle^2$$

$$= \sum_i x_i y_i + 0.5 \sum_{ij} \phi_{ij}(\mathbf{x}) \phi_{ij}(\mathbf{y})$$

$$= \langle (x_1, \dots, x_m, \sqrt{0.5}x_1x_1, \dots, \sqrt{0.5}x_mx_m), \\ (y_1, \dots, y_m, \sqrt{0.5}y_1y_1, \dots, \sqrt{0.5}y_my_m) \rangle$$

The Kernel Trick



What about:

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle + 0.5 \langle \mathbf{x}, \mathbf{y} \rangle^2$$

$$= \sum_i x_i y_i + 0.5 \sum_{ij} \phi_{ij}(\mathbf{x}) \phi_{ij}(\mathbf{y})$$

$$= \langle (x_1, \dots, x_m, \sqrt{0.5}x_1x_1, \dots, \sqrt{0.5}x_mx_m), \\ (y_1, \dots, y_m, \sqrt{0.5}y_1y_1, \dots, \sqrt{0.5}y_my_m) \rangle$$

The Kernel Trick

What about:

$$K(\mathbf{x}, \mathbf{y}) = 1 + \langle \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \langle \mathbf{x}, \mathbf{y} \rangle^2 + \frac{1}{6} \langle \mathbf{x}, \mathbf{y} \rangle^3 + \frac{1}{24} \langle \mathbf{x}, \mathbf{y} \rangle^4?$$

The Kernel Trick

What about:

$$K(x, y) = \sum_{i=0}^{\infty} \frac{\langle x, y \rangle^i}{i!}$$

The Kernel Trick

What about:

$$K(x, y) = \sum_{i=0}^{\infty} \frac{\langle x, y \rangle^i}{i!} = \exp \langle x, y \rangle$$

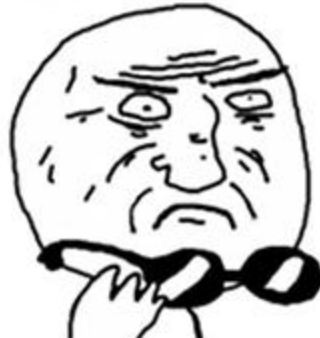
The Kernel Trick

What about:

$$K(x, y) = \sum_{i=0}^{\infty} \frac{\langle x, y \rangle^i}{i!} = \exp \langle x, y \rangle$$

Infinite-dimensional feature space!

MOTHER OF GOD...



The Kernel Trick

Gaussian (RBF) kernel

$$\begin{aligned} K(x, y) &= \\ &= \exp(-\gamma \|x - y\|^2) \\ &= \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right) \end{aligned}$$



The Kernel Trick

Exponential kernel

$$K(x, y) = \exp\left(-\frac{\|x - y\|}{2\sigma^2}\right)$$



Kernels

$$k(x, y) = x^T y + c \quad k(x, y) = \sqrt{\|x - y\|^2 + c^2} \quad k(x, y) = \frac{\theta}{\|x - y\|} \sin \frac{\|x - y\|}{\theta}$$

$$k(x, y) = (\alpha x^T y + c)^d \quad k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right) \quad k(x, y) = \frac{1}{\sqrt{\|x - y\|^2 + c^2}}$$

$$k(x, y) = \exp\left(-\frac{\|x - y\|}{2\sigma^2}\right) \quad k(x, y) = \frac{2}{\pi} \arccos\left(-\frac{\|x - y\|}{\sigma}\right) - \frac{2}{\pi} \frac{\|x - y\|}{\sigma} \sqrt{1 - \left(\frac{\|x - y\|}{\sigma}\right)^2}$$

$$k(x, y) = \sum_{k=1}^n \exp(-\sigma(x^k - y^k)^2) \quad k(x, y) = 1 - \frac{3}{2} \frac{\|x - y\|}{\sigma} + \frac{1}{2} \left(\frac{\|x - y\|}{\sigma}\right)^3 \quad k(x, y) = -\log(\|x - y\|^d + 1)$$

$$k(x, y) = \tanh(\alpha x^T y + c) \quad k(x, y) = -\|x - y\|^d \quad k(x, y) = \frac{1}{1 + \frac{\|x - y\|^2}{\sigma}}$$

$$k(x, y) = 1 - \frac{\|x - y\|^2}{\|x - y\|^2 + c} \quad k(x, y) = 1 + xy + xy \min(x, y) - \frac{x + y}{2} \min(x, y)^2 + \frac{1}{3} \min(x, y)^3$$

$$k(x, y) = B_{2p+1}(x - y) \quad k(x, y) = \sum_{i=1}^n \min(x_i, y_i) \quad k(x, y) = \prod_{i=1}^N h\left(\frac{x_i - c}{a}\right) h\left(\frac{y_i - c}{a}\right)$$

$$k(x, y) = \sum_{i=1}^m \min(|x_i|^\alpha, |y_i|^\beta) \quad k(x, y) = \frac{J_{v+1}(\sigma\|x - y\|)}{\|x - y\|^{-n(v+1)}} \quad k(x, y) = 1 - \sum_{i=1}^n \frac{(x_i - y_i)^2}{\frac{1}{2}(x_i + y_i)}$$

$$\kappa_l(a, b) = \sum_{c \in \{0,1\}} P(Y = c \mid X_l = a) P(Y = c \mid X_l = b) \quad k(x, y) = \frac{1}{1 + \|x - y\|^d}$$



Structured data kernels

▶ String kernels

- ▶ P-spectrum kernels
- ▶ All-subsequences kernels
- ▶ Gap-weighted subsequences kernels
- ▶ ...

▶ Graph & tree kernels

- ▶ Co-rooted subtrees
- ▶ All subtrees
- ▶ Random walks
- ▶ ...



Kernel

- ▶ A function $K(\mathbf{x}, \mathbf{y})$ is a *kernel*, if

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

for some *feature map* ϕ .

Kernel matrix

- ▶ For a given kernel function K and a finite dataset $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ the $n \times n$ matrix

$$\mathbf{K}_{ij} := K(\mathbf{x}_i, \mathbf{x}_j)$$

is called the *kernel matrix*.

Kernel matrix

- ▶ Let X be the data matrix, then

$$K = XX^T$$

is the kernel matrix for the *linear* kernel

$$K(x, y) = x^T y$$

Kernel matrix

- ▶ Let X be the data matrix, then

$$K = XX^T$$

is the kernel matrix for the *linear* kernel

$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$$

- ▶ Let ϕ be a feature mapping. Then*

$$K = \phi(X)\phi(X)^T$$

is the kernel matrix for the corresponding kernel $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$.

Kernel theorem

- ▶ Not every function K is a kernel!
e. g. $K(x, y) = -1$ is not
- ▶ Not every $n \times n$ matrix is a Kernel matrix!

Kernel theorem

- ▶ Theorem:

K is a kernel function $\Leftrightarrow K$ is *symmetric positive semidefinite*

- ▶ A function is *positive semidefinite* iff for any finite dataset $\{x_1, x_2, \dots, x_n\}$ the corresponding *kernel matrix* is positive semidefinite.

Kernel closure

- ★ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$
- ★ $\kappa(x, z) = \alpha \kappa_1(x, z)$
- ★ $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$
- ★ $\kappa(x, z) = f(x)f(z)$ where f is a real-valued function
- ★ $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$
- ★ $\kappa(x, z) = x^T B z$ where B is a psd matrix.

P. Agius – L3, Spring 2008

Kernel closure

Feature space concatenation

- ★ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$
- ★ $\kappa(x, z) = \alpha \kappa_1(x, z)$
- ★ $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$
- ★ $\kappa(x, z) = f(x)f(z)$ where f is a real-valued function
- ★ $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$
- ★ $\kappa(x, z) = x^T B z$ where B is a psd matrix.

P. Agius – L3, Spring 2008

Kernel closure

Feature space scaling

- ★ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$
- ★ $\kappa(x, z) = \alpha \kappa_1(x, z)$
- ★ $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$
- ★ $\kappa(x, z) = f(x)f(z)$ where f is a real-valued function
- ★ $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$
- ★ $\kappa(x, z) = x^T B z$ where B is a psd matrix.

P. Agius – L3, Spring 2008

Kernel closure

- ★ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$
- ★ $\kappa(x, z) = \alpha \kappa_1(x, z)$
- ★ $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$
- ★ $\kappa(x, z) = f(x)f(z)$ where f is a real-valued function
- ★ $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$
- ★ $\kappa(x, z) = x^T B z$ where B is a psd matrix.

Feature space tensor product

P. Agius – L3, Spring 2008

Kernel closure

- ★ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$
- ★ $\kappa(x, z) = \alpha \kappa_1(x, z)$
- ★ $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$
- ★ $\kappa(x, z) = f(x)f(z)$ where f is a real-valued function
- ★ $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$
- ★ $\kappa(x, z) = x^T B z$ where B is a psd matrix.

Feature map composition

P. Agius – L3, Spring 2008

The Kernel-based Approach

Nonlinear feature mapping

The Kernel trick

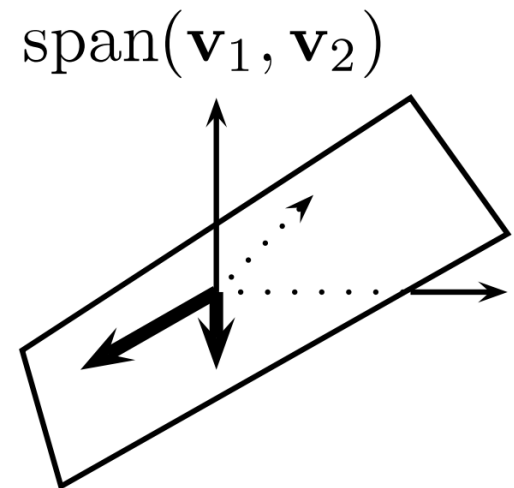
Dual representation



Dual representation

Consider the *linear span* of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, i.e. the set of all vectors \mathbf{w} of the form

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i = \mathbf{X}^T \boldsymbol{\alpha}$$



Dual representation

Whenever

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i = \mathbf{X}^T \boldsymbol{\alpha}$$

we shall refer to $\boldsymbol{\alpha}$ as the **dual representation** of \mathbf{w} .

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} =$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} =$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T (\boldsymbol{\alpha} + \boldsymbol{\beta})$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T(2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T(\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$\langle \mathbf{w}, \mathbf{u} \rangle =$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T (\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$\langle \mathbf{w}, \mathbf{u} \rangle = \mathbf{w}^T \mathbf{u} = \boldsymbol{\alpha}^T \mathbf{X} \mathbf{X}^T \boldsymbol{\beta} = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\beta}$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T (\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$\langle \mathbf{w}, \mathbf{u} \rangle = \mathbf{w}^T \mathbf{u} = \boldsymbol{\alpha}^T \mathbf{X} \mathbf{X}^T \boldsymbol{\beta} = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\beta}$$

$$\|\mathbf{w} - \mathbf{u}\|^2 =$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T (2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T (\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$\langle \mathbf{w}, \mathbf{u} \rangle = \mathbf{w}^T \mathbf{u} = \boldsymbol{\alpha}^T \mathbf{X} \mathbf{X}^T \boldsymbol{\beta} = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\beta}$$

$$\|\mathbf{w} - \mathbf{u}\|^2 = \langle \mathbf{w} - \mathbf{u}, \mathbf{w} - \mathbf{u} \rangle = \dots$$

Dual coordinates

Let

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{u} = \mathbf{X}^T \boldsymbol{\beta}$$

Then

$$2\mathbf{w} = \mathbf{X}^T(2\boldsymbol{\alpha})$$

$$\mathbf{w} + \mathbf{u} = \mathbf{X}^T(\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$\langle \mathbf{w}, \mathbf{u} \rangle = \mathbf{w}^T \mathbf{u} = \boldsymbol{\alpha}^T \mathbf{X} \mathbf{X}^T \boldsymbol{\beta} = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\beta}$$

$$\|\mathbf{w} - \mathbf{u}\|^2 = \langle \mathbf{w} - \mathbf{u}, \mathbf{w} - \mathbf{u} \rangle = \dots$$



So what?



The Representer Theorem

The SVM weight vector lies in the span of the data points, i.e.

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i = \mathbf{X}^T \boldsymbol{\alpha}$$

for some $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$

The Representer Theorem

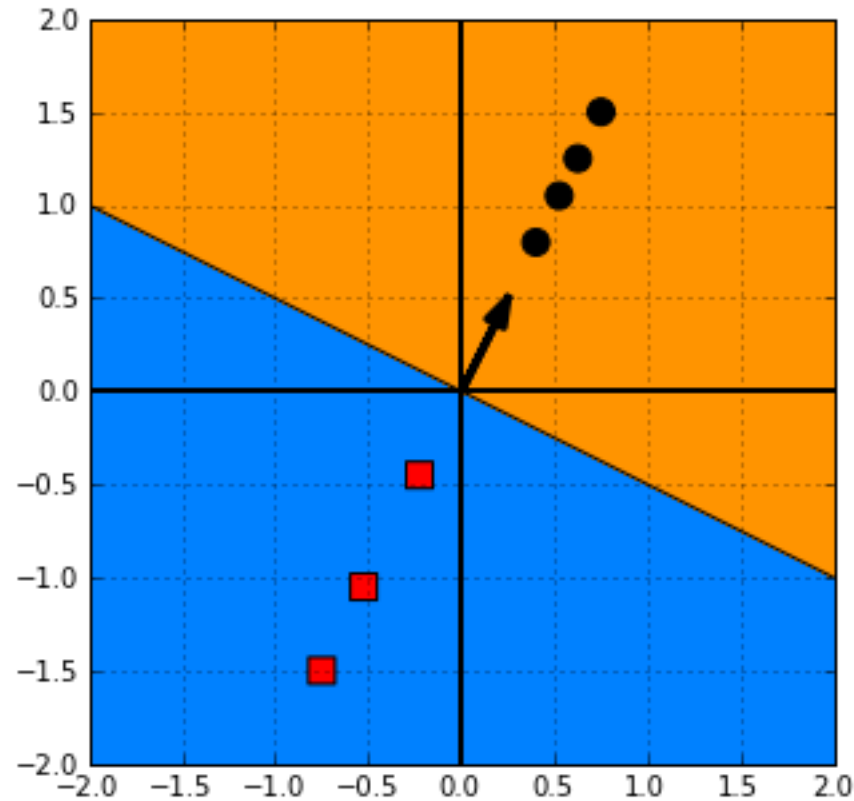
The SVM weight vector lies in the span of the data points, i.e.

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i = \mathbf{X}^T \boldsymbol{\alpha}$$

for some $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$

- Easy to prove
- Actually holds for **pretty much any linear model** with ℓ_2 -penalty.

The Representer Theorem



Example

- ▶ Recall OLS regression:

$$\mathbf{w} = \underline{\hspace{2cm}}$$

Example

- ▶ Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

Example

- Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$
$$\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}$$

Example

- Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T [(\mathbf{X}\mathbf{X}^T)^+ \mathbf{y}]$$

Example

- Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T [(\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}]$$

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

Kernelization

1. $x_i \rightarrow \phi(x_i)$

Nonlinear feature mapping

Kernelization

1. $\mathbf{x}_i \rightarrow \phi(\mathbf{x}_i)$

Nonlinear feature mapping

2. $\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i)$

Dual representation

Kernelization

1. $\mathbf{x}_i \rightarrow \phi(\mathbf{x}_i)$

Nonlinear feature mapping

2. $\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i)$

Dual representation

3.

$$\begin{aligned} f(\mathbf{z}) &= \langle \mathbf{w}, \phi(\mathbf{z}) \rangle + w_0 = \left\langle \sum_i \alpha_i \phi(\mathbf{x}_i), \phi(\mathbf{z}) \right\rangle + w_0 = \\ &= \sum_i \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{z}) \rangle + w_0 \\ &= \sum_i \alpha_i K(\mathbf{x}_i, \mathbf{z}) + w_0 \end{aligned}$$

The Kernel trick

Example

- Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T [(\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}]$$

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

Example

- Recall OLS regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^T [(\mathbf{X} \mathbf{X}^T)^+ \mathbf{y}]$$

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha}$$

$$\boldsymbol{\alpha} = \mathbf{K}^+ \mathbf{y} = \mathbf{K}^{-1} \mathbf{y}$$

Kernelization: Summary

- ▶ Take a linear learning algorithm

$$\mathbf{X}, \mathbf{y} \rightarrow \mathbf{w}, w_0$$

- ▶ Rewrite it to use only inner products of data points and return the dual representation of \mathbf{w}

$$\mathbf{K}, \mathbf{y} \rightarrow \boldsymbol{\alpha}, w_0$$

- ▶ Plug in any kernel and play!
- ▶ Congratulations, you've got a nonlinear model!

Quiz

The three ingredients of kernel methods:

- ▶ _____
- ▶ _____
- ▶ _____

- ▶ Function/matrix K is a kernel function/matrix iff it is _____
- ▶ Dual representation: _____ = _____

Quiz



Those algorithms have kernelized versions:

_____ . . .



Questions?

