

---

# Machine Learning

## Math Essentials Part 2

# Gaussian distribution

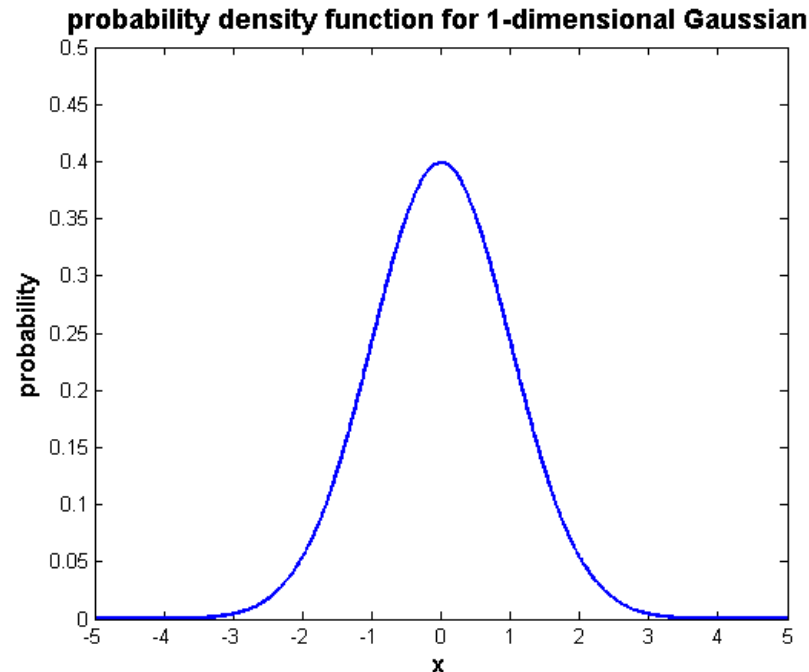
---

- Most commonly used continuous probability distribution
- Also known as the normal distribution
- Two parameters define a Gaussian:
  - Mean  $\mu$  location of center
  - Variance  $\sigma^2$  width of curve

# Gaussian distribution

## In one dimension

$$N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



# Gaussian distribution

## In one dimension

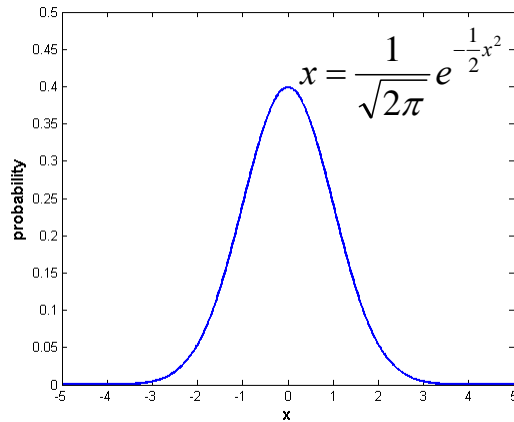
$$N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Causes pdf to decrease as distance from center increases

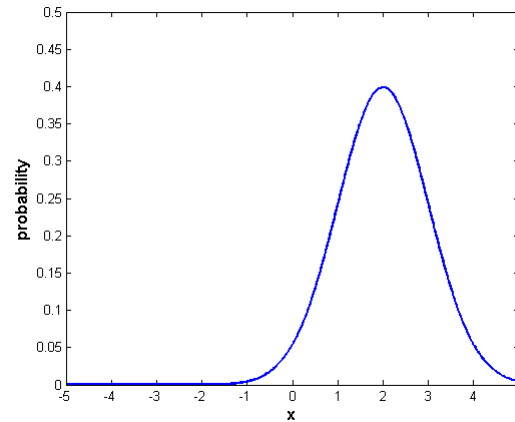
Controls width of curve

Normalizing constant: insures that distribution integrates to 1

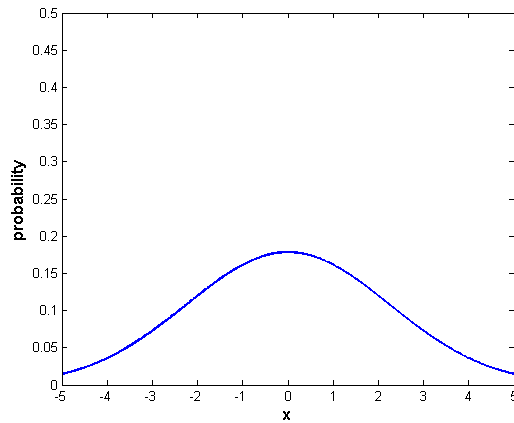
# Gaussian distribution



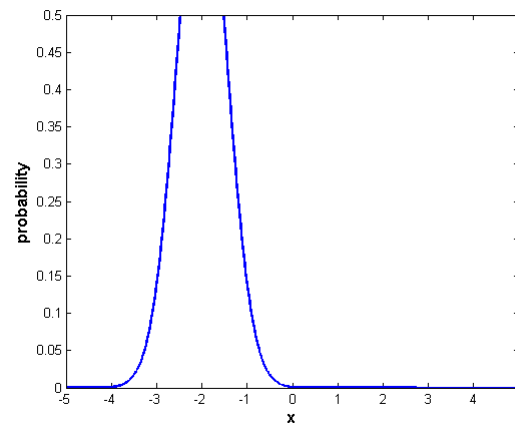
$$\mu = 0 \quad \sigma^2 = 1$$



$$\mu = 2 \quad \sigma^2 = 1$$



$$\mu = 0 \quad \sigma^2 = 5$$



$$\mu = -2 \quad \sigma^2 = 0.3$$

# Multivariate Gaussian distribution

---

## In $d$ dimensions

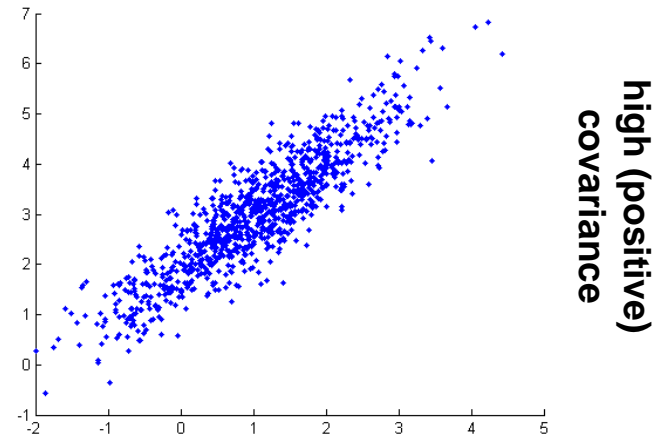
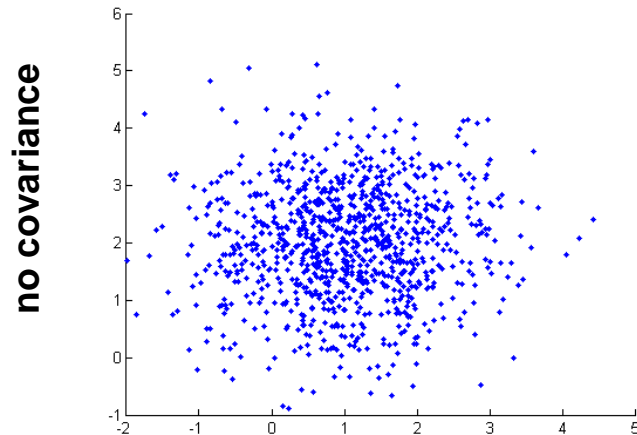
$$N(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

- $\mathbf{x}$  and  $\boldsymbol{\mu}$  now  $d$ -dimensional vectors
  - $\boldsymbol{\mu}$  gives center of distribution in  $d$ -dimensional space
- $\sigma^2$  replaced by  $\boldsymbol{\Sigma}$ , the  $d \times d$  covariance matrix
  - $\boldsymbol{\Sigma}$  contains pairwise covariances of every pair of features
  - Diagonal elements of  $\boldsymbol{\Sigma}$  are variances  $\sigma^2$  of individual features
  - $\boldsymbol{\Sigma}$  describes distribution's shape and spread

# Multivariate Gaussian distribution

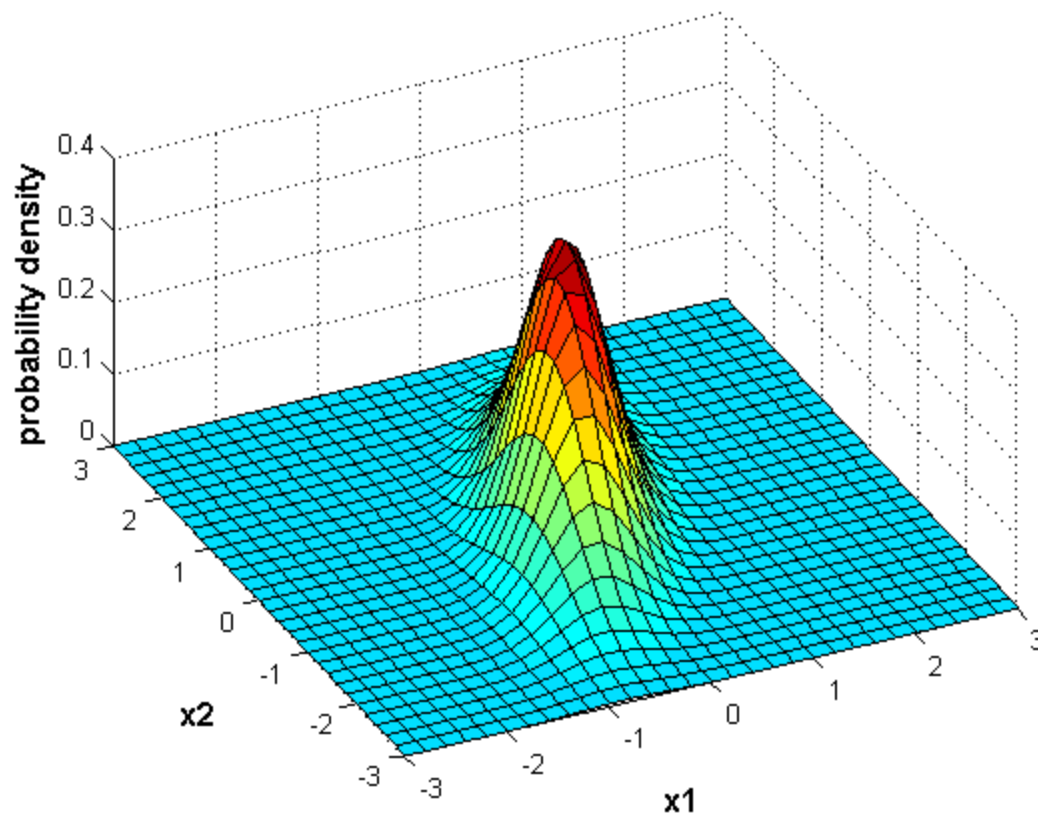
- Covariance

- Measures tendency for two variables to deviate from their means in same (or opposite) directions at same time



# Multivariate Gaussian distribution

## In two dimensions

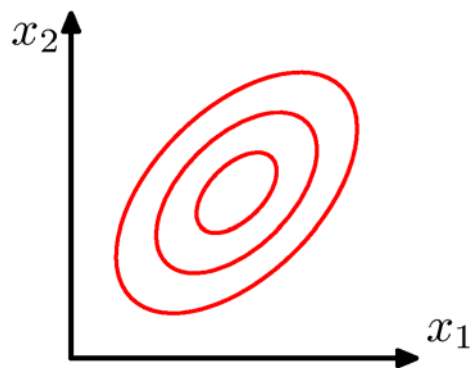


$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$$\boldsymbol{\Sigma} = \begin{bmatrix} 0.25 & 0.3 \\ 0.3 & 1 \end{bmatrix}$$



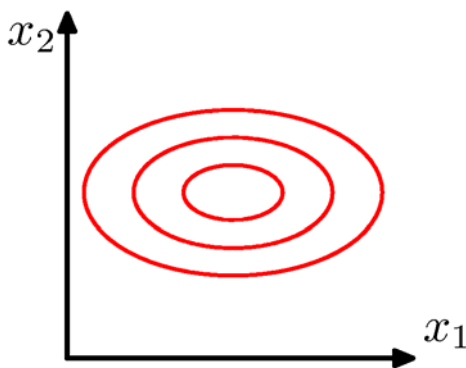
# Multivariate Gaussian distribution

## In two dimensions



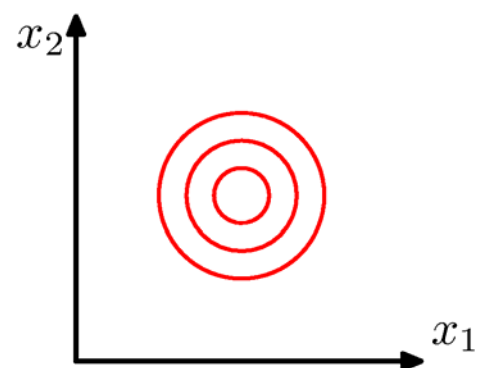
(a)

$$\Sigma = \begin{bmatrix} 2 & 0.6 \\ 0.6 & 2 \end{bmatrix}$$



(b)

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

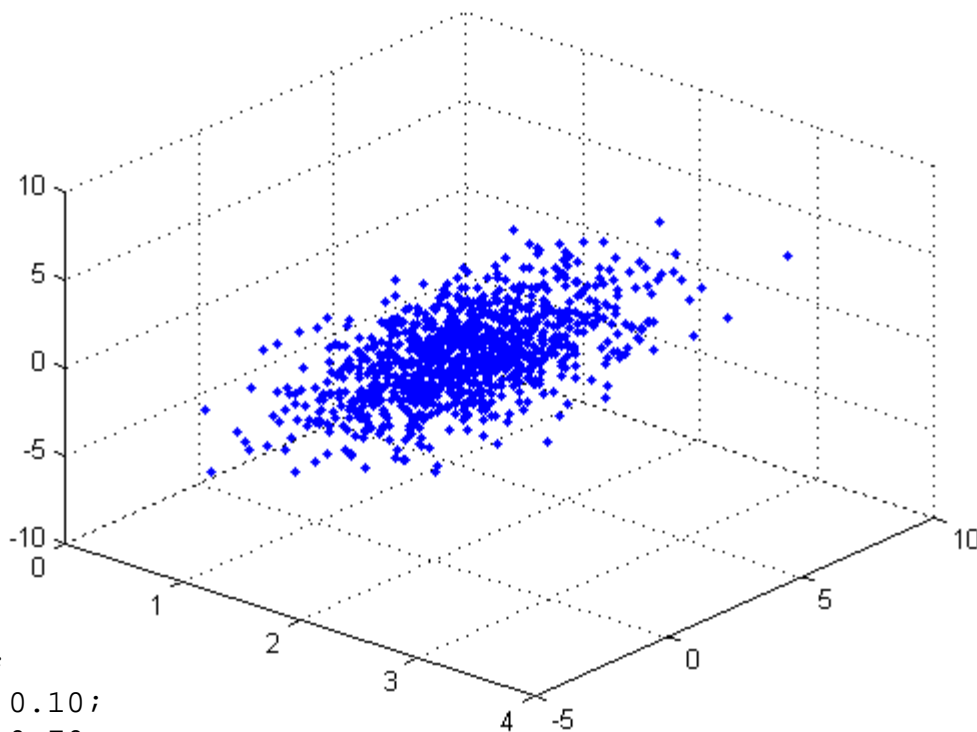


(c)

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

# Multivariate Gaussian distribution

## In three dimensions



```
rng( 1 );  
mu = [ 2; 1; 1 ];  
sigma = [ 0.25 0.30 0.10;  
          0.30 1.00 0.70;  
          0.10 0.70 2.00] ;  
x = randn( 1000, 3 );  
x = x * sigma;  
x = x + repmat( mu', 1000, 1 );  
scatter3( x( :, 1 ), x( :, 2 ), x( :, 3 ), '.' );
```

# Vector projection

- Orthogonal projection of  $\mathbf{y}$  onto  $\mathbf{x}$ 
  - Can take place in any space of dimensionality  $\geq 2$

- Unit vector in direction of  $\mathbf{x}$  is

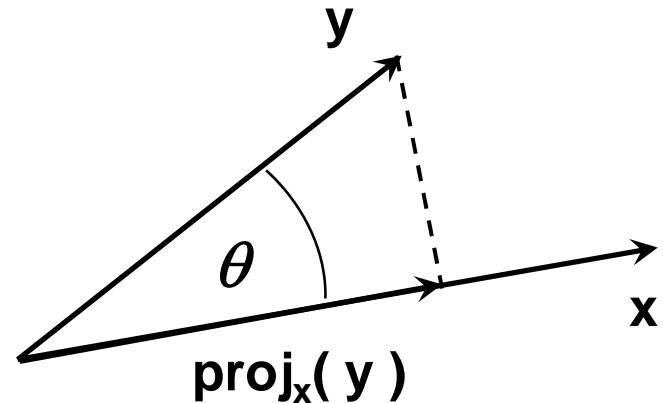
$$\mathbf{x} / \|\mathbf{x}\|$$

- Length of projection of  $\mathbf{y}$  in direction of  $\mathbf{x}$  is

$$\|\mathbf{y}\| \cdot \cos(\theta)$$

- Orthogonal projection of  $\mathbf{y}$  onto  $\mathbf{x}$  is the vector

$$\mathbf{proj}_x(\mathbf{y}) = \mathbf{x} \cdot \|\mathbf{y}\| \cdot \cos(\theta) / \|\mathbf{x}\| = [(\mathbf{x} \cdot \mathbf{y}) / \|\mathbf{x}\|^2] \mathbf{x} \quad (\text{using dot product alternate form})$$



# Linear models

- There are many types of **linear** models in machine learning.
  - Common in both classification and regression.
  - A linear model consists of a vector  $\mathbf{w}$  in  $d$ -dimensional feature space.
  - The vector  $\mathbf{w}$  attempts to capture the strongest gradient (rate of change) in the output variable, as seen across all training samples.
  - Different linear models optimize  $\mathbf{w}$  in different ways.
  - A point  $\mathbf{x}$  in feature space is mapped from  $d$  dimensions to a scalar (1-dimensional) output  $z$  by projection onto  $\mathbf{w}$ :

$$z = w_0 + \mathbf{W} \cdot \mathbf{X} = w_0 + w_1 x_1 + \cdots + w_d x_d$$

cf. Lecture 5b

$w_0 \equiv \alpha$

$\mathbf{w} \equiv \beta$

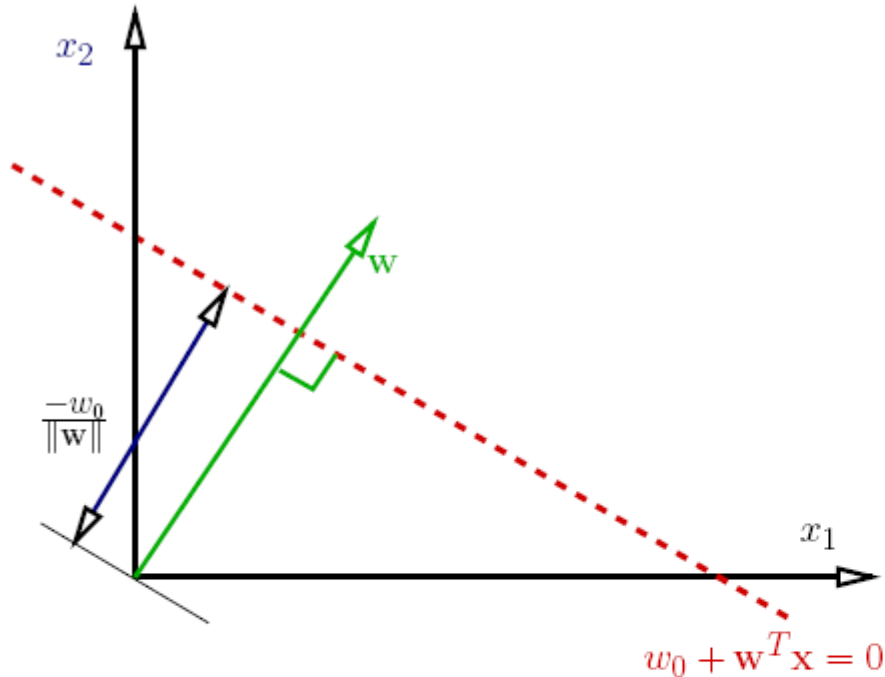
# Linear models

- There are many types of **linear** models in machine learning.
  - The projection output  $z$  is typically transformed to a final predicted output  $y$  by some function  $f$ :

$$y = f(z) = f(w_0 + \mathbf{w} \cdot \mathbf{x}) = f(w_0 + w_1x_1 + \dots + w_dx_d)$$

- ◆ example: for logistic regression,  $f$  is logistic function
- ◆ example: for linear regression,  $f(z) = z$
- Models are called linear because they are a linear function of the model vector components  $w_1, \dots, w_d$ .
- Key feature of all linear models: no matter what  $f$  is, a constant value of  $z$  is transformed to a constant value of  $y$ , so decision boundaries remain linear even after transform.

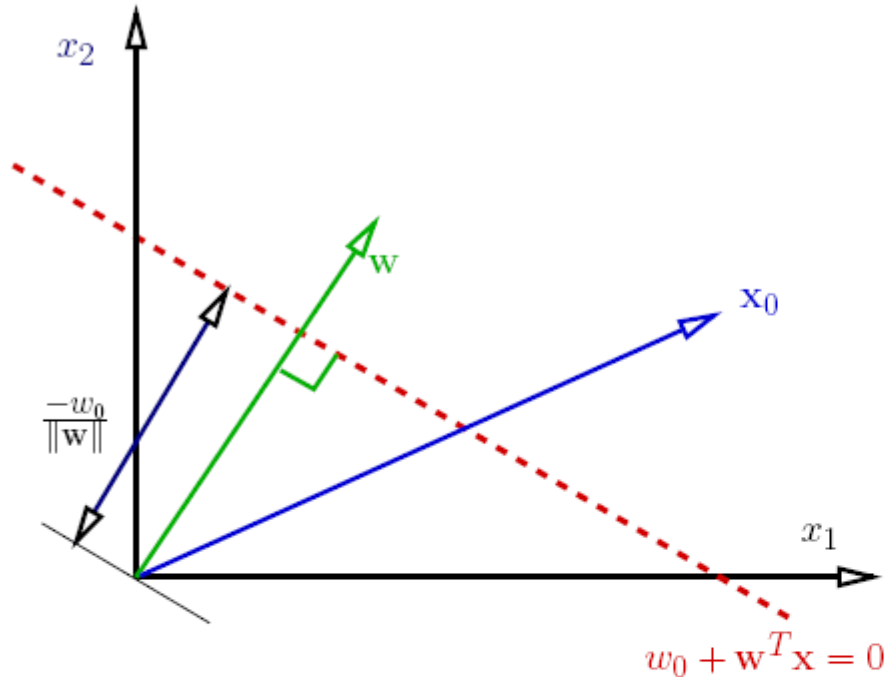
# Geometry of projections



- $\mathbf{w}^T \mathbf{x} = 0$ : a line passing through the origin and *orthogonal* to  $\mathbf{w}$
- $\mathbf{w}^T \mathbf{x} + w_0 = 0$  shifts the line along  $\mathbf{w}$ .

slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)

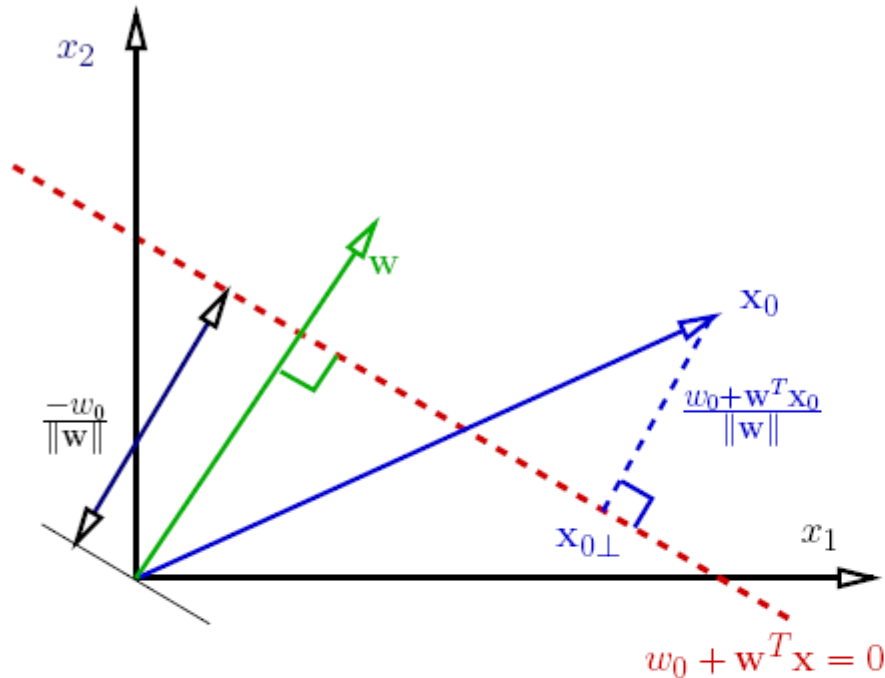
# Geometry of projections



- $w^T x = 0$ : a line passing through the origin and *orthogonal* to  $w$
- $w^T x + w_0 = 0$  shifts the line along  $w$ .

slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)

# Geometry of projections

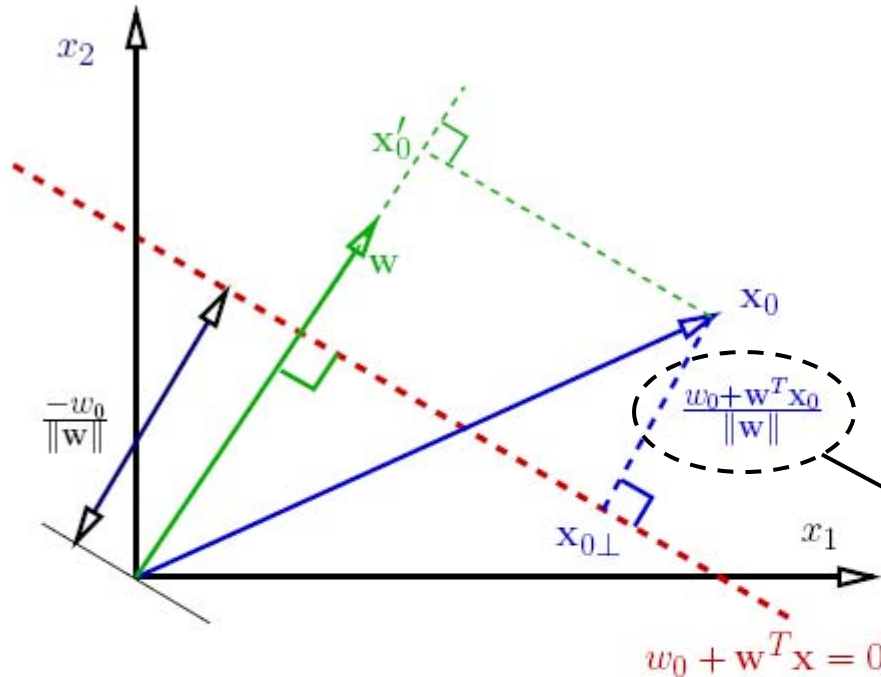


- $w^T x = 0$ : a line passing through the origin and *orthogonal* to  $w$
- $w^T x + w_0 = 0$  shifts the line along  $w$ .

slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)



# Geometry of projections



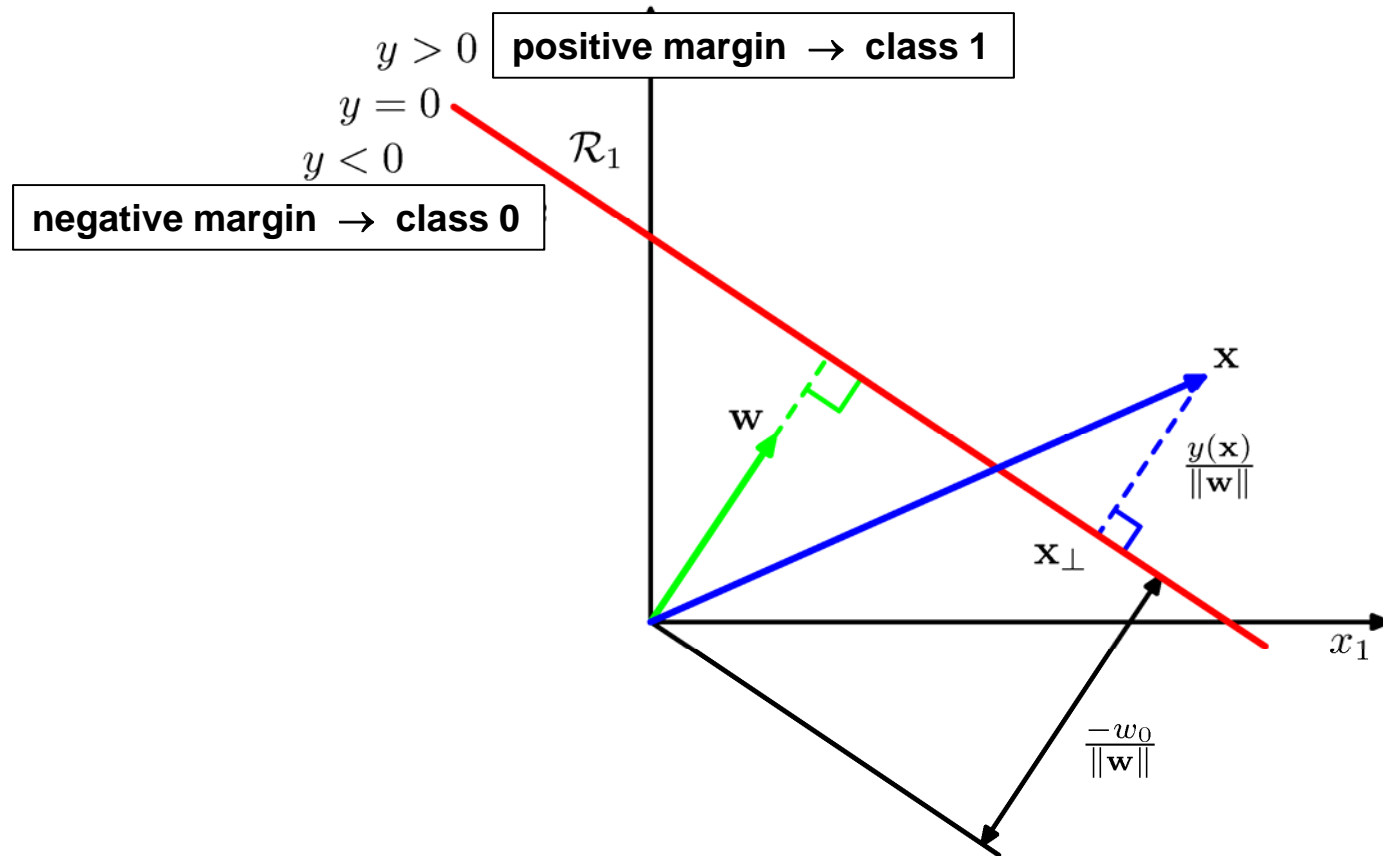
- $w^T \mathbf{x} = 0$ : a line passing through the origin and *orthogonal* to  $w$
- $w^T \mathbf{x} + w_0 = 0$  shifts the line along  $w$ .

margin

- $x'$  is the projection of  $x$  on  $w$ .
- Set up a new 1D coordinate system:  $x \rightarrow \left( \frac{w_0 + x^T \mathbf{x}}{\|\mathbf{w}\|} \right)$ .

slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)

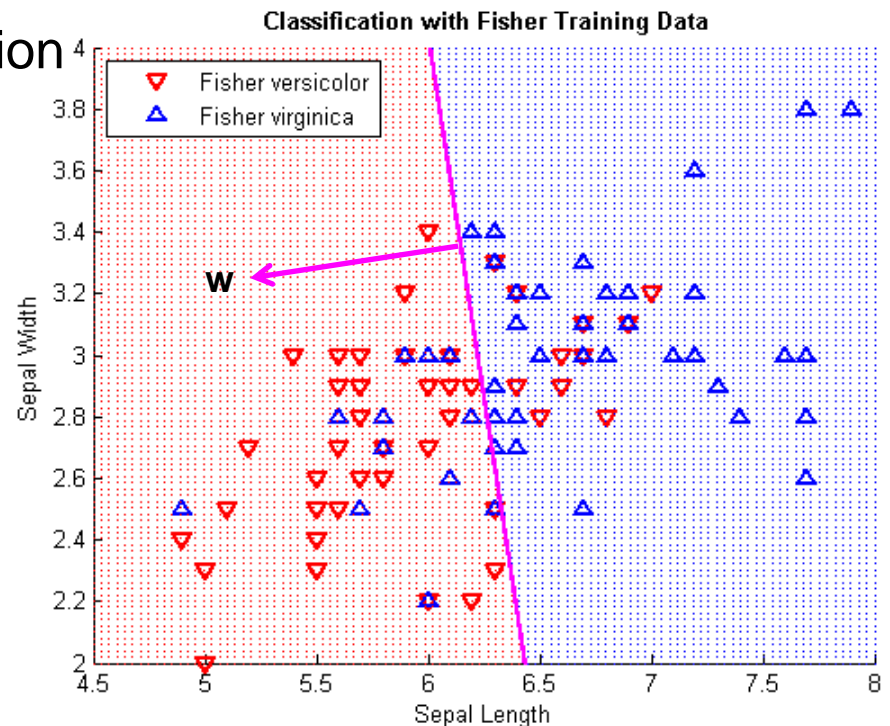
# From projection to prediction



# Logistic regression in two dimensions

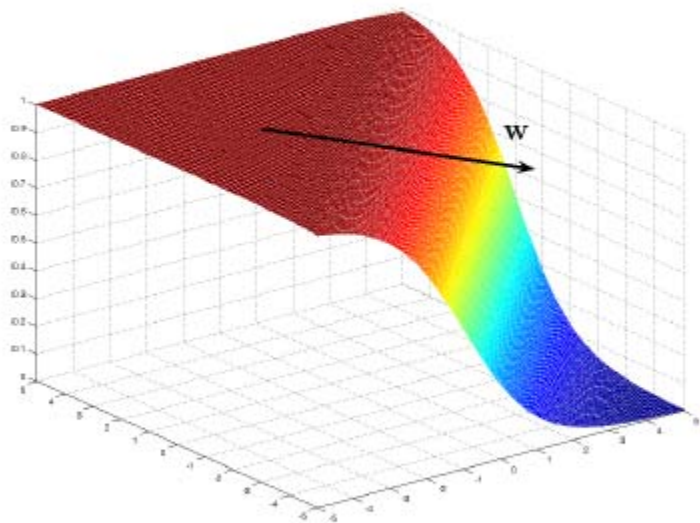
## Interpreting the model vector of coefficients

- From MATLAB:  $B = [ 13.0460 \quad -1.9024 \quad -0.4047 ]$
- $w_0 = B(1)$ ,  $\mathbf{w} = [ w_1 \ w_2 ] = B(2:3)$
- $w_0$ ,  $\mathbf{w}$  define location and orientation of decision boundary
  - $w_0$  is distance of decision boundary from origin
  - decision boundary is perpendicular to  $\mathbf{w}$
- magnitude of  $\mathbf{w}$  defines gradient of probabilities between 0 and 1



# Logistic function in $d$ dimensions

- What if  $\mathbf{x} \in \mathbb{R}^d = [x_1 \dots x_d]^T$ ?
- $\sigma(w_0 + \mathbf{w}^T \mathbf{x})$  is a scalar function of a scalar variable  $w_0 + \mathbf{w}^T \mathbf{x}$ .



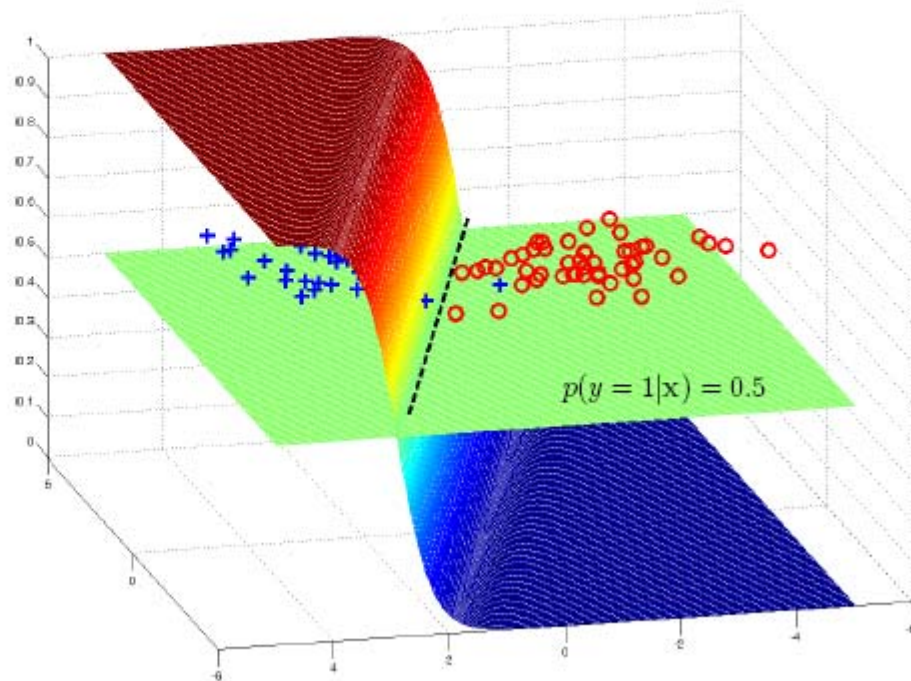
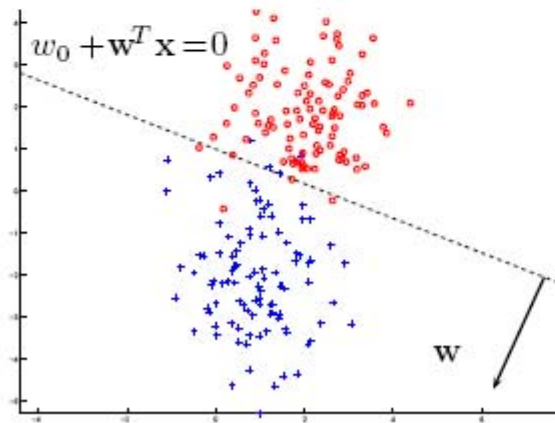
- the direction of  $\mathbf{w}$  determines orientation;
- $w_0$  determines the location;
- $\|\mathbf{w}\|$  determines the slope.

slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)

# Decision boundary for logistic regression

$$p(y = 1 | \mathbf{x}) = \sigma(w_0 + \mathbf{w}^T \mathbf{x}) = 1/2 \Leftrightarrow w_0 + \mathbf{w}^T \mathbf{x} = 0$$

- With linear logistic model we get a linear decision boundary.



slide thanks to Greg Shakhnarovich (CS195-5, Brown Univ., 2006)