# Big Picture

Machine Learning – 10701/15781
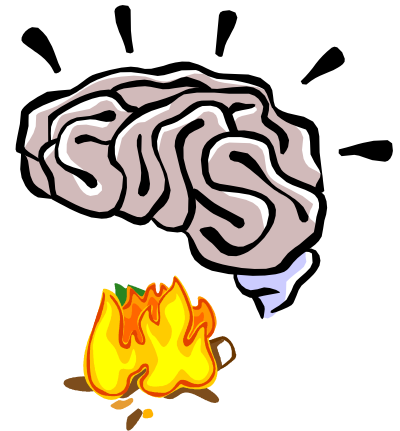
Carlos Guestrin

Carnegie Mellon University

March 2nd, 2005

# What you have learned thus far

- Learning is function approximation
- Point estimation
- Regression
- Naïve Bayes
- Logistic regression
- Bias-Variance tradeoff
- Neural nets
- Decision trees
- Cross validation
- Boosting
- Instance-based learning
- SVMs
- Kernel trick
- PAC learning
- VC dimension
- Margin bounds
- Mistake bounds

# Review material in terms of…

- Types of learning problems

- Hypothesis spaces

- Loss functions

- Optimization algorithms

# Text Classification

text $\longrightarrow$ {C, P, U, ...}
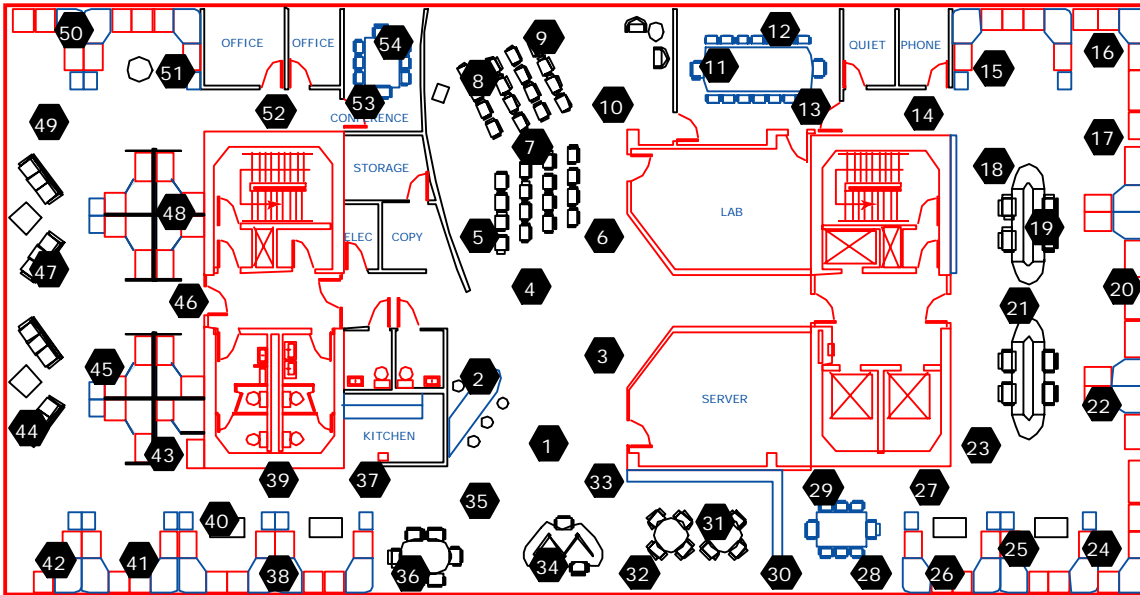


$\longrightarrow$ Company home page
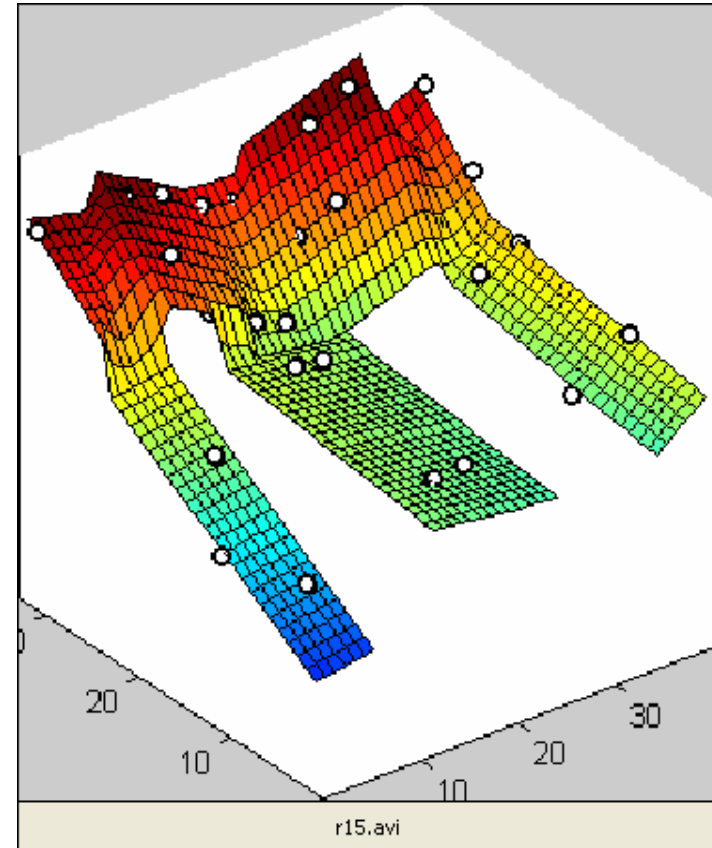
vs

Personal home page

vs

Univeristy home page
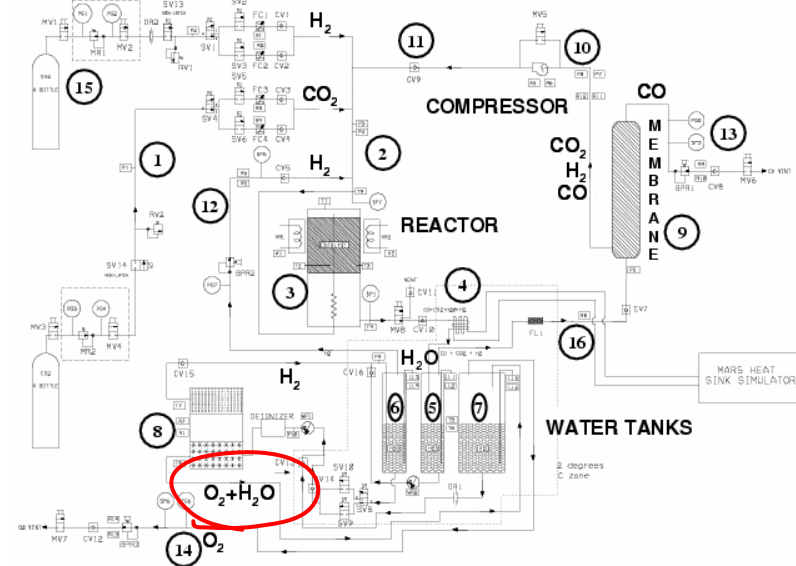
vs

…

# Function fitting

$x, y, t \rightarrow temps$



Temperature data



r15.avi

# Monitoring a complex system



- Reverse water gas shift system (RWGS)
- Learn model of system from data
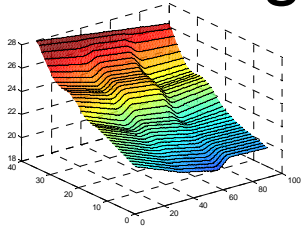- Use model to predict behavior and detect faults

# Types of learning problems

- Classification

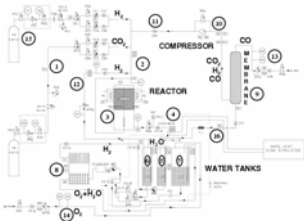text $\longrightarrow \{C, P, U, \ldots\}$

- Regression

$x, y, t \longrightarrow \mathbb{R}$

- Density estimation

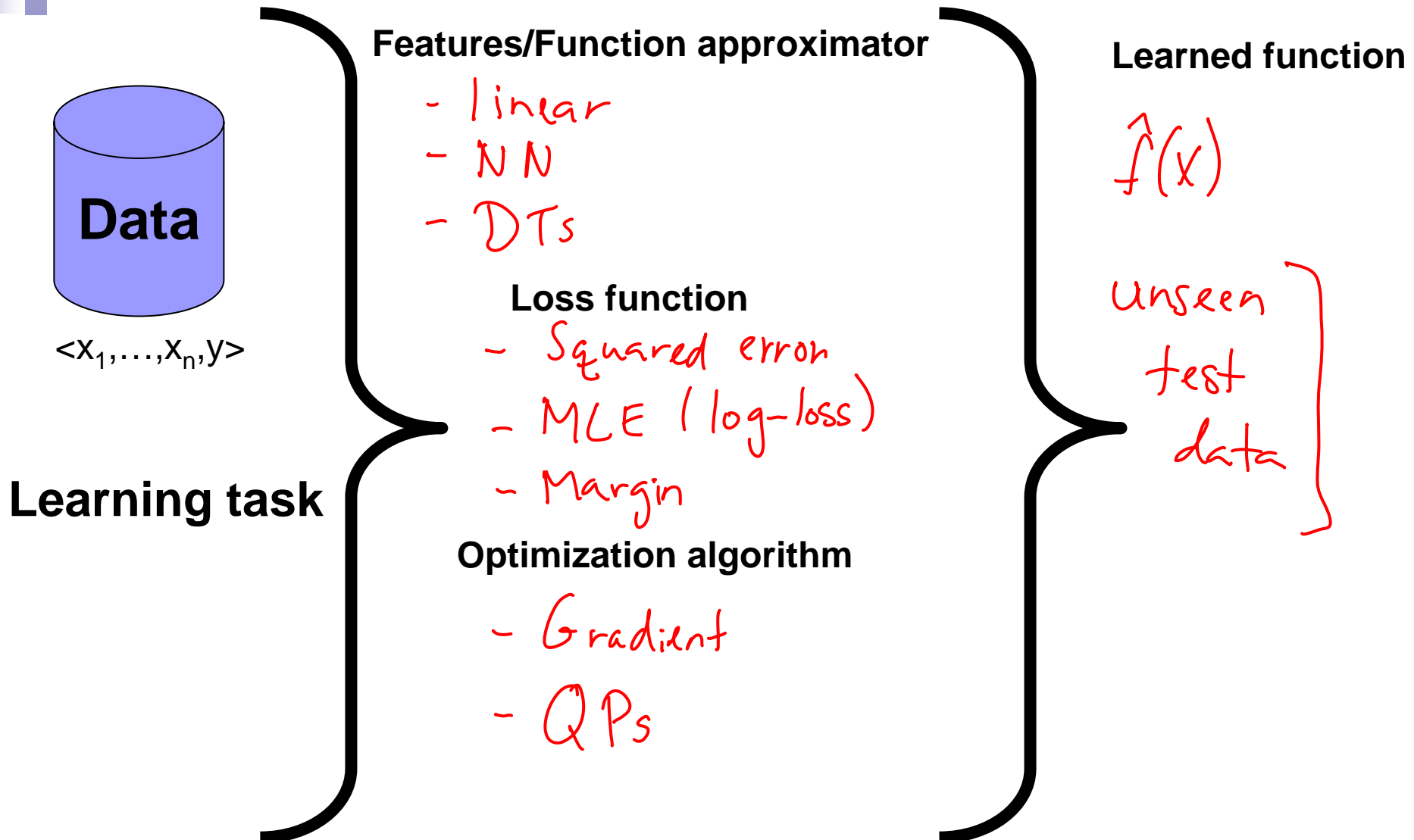sensors $\longrightarrow [0, 1]$

**Input – Features**

$X \qquad \phi(X)$

**Output?**

$Y$: classification; discrete

Regression: $\mathbb{R}$

Density Est.: $[0, 1]$

# The learning problem

**Data**

$<x_1,\ldots,x_n,y>$

**Learning task**

**Features/Function approximator**
- linear
- NN
- DTs

**Loss function**
- Squared error
- MLE (log-loss)
- Margin

**Optimization algorithm**
- Gradient
- QPs

**Learned function**

$\hat{f}(x)$

Unseen test data

# Comparing learning algorithms

- Hypothesis space

- Loss function

- Optimization algorithm

# Naïve Bayes versus Logistic regression

Density estimation: $P(Y|x)$

**Naïve Bayes**

**Logistic regression**

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(Y=1|x) = \frac{1}{1 + exp(w_0 + \sum_i w_i x_i)}$$

$$P(X|Y) = \prod_i P(X_i|Y)$$

weaker indep. assumptions

Strong indep. assumption

loss: max $P(Y|x)$

loss function:
max $P(X,Y)$

Gradient, MLE

MLE

equivalence

indep. Gaussian features

# Naïve Bayes versus Logistic regression – Classification as density estimation

$$P(Y|X)$$

- Choose class with highest probability

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

- In addition to class, we get certainty measure

# Logistic regression versus Boosting

## Logistic regression

$$P(Y = y_i | \mathbf{x}) = \frac{1}{1 + exp(-y_i(\mathbf{w}.\mathbf{x} + b))}$$

Log-loss

$$\sum_{j=1}^{m} \log \left[ 1 + exp(-y_i(\mathbf{w}.\mathbf{x}_j + b)) \right]$$

$\rightarrow$ Sign $(w.x+b)$ classifier

$\rightarrow$ select features apriori

Global optima

## Boosting

optimizing $\alpha_t$

Classifier

$$sign \left( \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x}) \right)$$

Exponential-loss

$$\frac{1}{m} \sum_{j=1}^{m} \exp \left( -y_j \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x_j}) \right)$$

$\neq$

$\rightarrow$ weak learners are the features $h_t(x)$

$\rightarrow$ sequential optimization

# Linear classifiers – Logistic regression versus SVMs

$w.x + b = 0$

# What's the difference between SVMs and Logistic Regression? (Revisited again)

| | **SVMs** | **Logistic Regression** |
|---|---|---|
| **Loss function** | Hinge loss | Log-loss |
| **High dimensional features with kernels** | Yes! | Yes! |
| **Solution sparse** | Often yes! | Almost always no! |
| | Classification | P(y |x) density estimation |

# SVMs and instance-based learning

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i)$$

$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

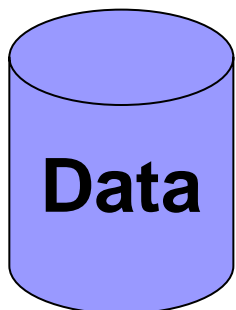for any $k$ where $C > \alpha_k > 0$

**SVMs**

**Classify as** → $sign\left(\sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) + b\right)$

## Instance based learning

as *density estimation*

*optimize* $\alpha_i$

"$\alpha_i$" are fixed

$$P(y \mid \mathbf{x}) = \frac{\sum_i y_i K(\mathbf{x}, \mathbf{x}_i)}{\sum_i K(\mathbf{x}, \mathbf{x}_i)} > 0.5?$$
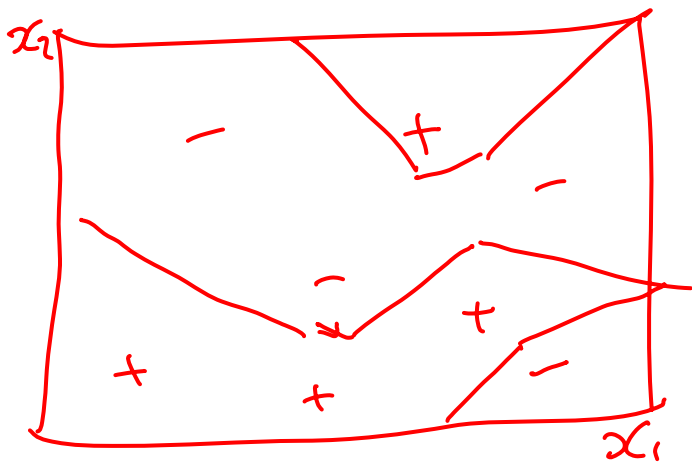
**Data**

**Classify as** →

$$sign\left(\sum_i y_i K(\mathbf{x}, \mathbf{x}_i) - 0.5 \sum_i K(\mathbf{x}, \mathbf{x}_i)\right)$$

$<x_1,\ldots,x_n,y>$

# Instance-based learning versus Decision trees

**1-Nearest neighbor**                    **Decision trees**
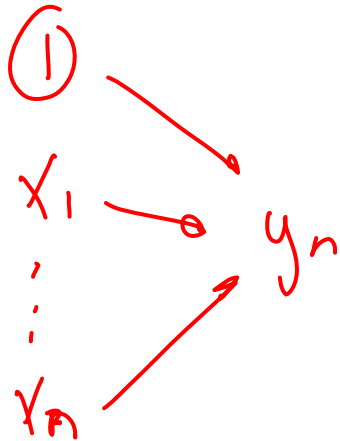


Voronoi split                            axis alligned split

# Logistic regression versus Neural nets

$$g\left(w_0 + \sum_i w_i x_i\right) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

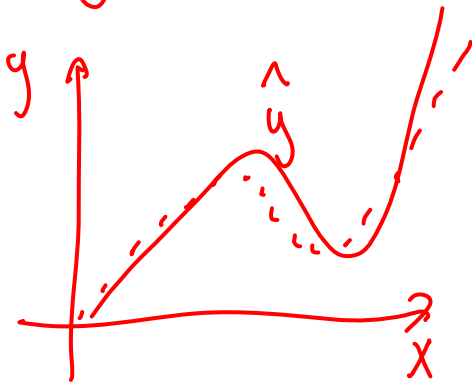**Logistic regression**  **Neural Nets**



loss:   log loss

loss:   Sum-squared error

# Linear regression versus Kernel regression

## Linear Regression

$$\hat{y} = w_0 + \sum_i w_i f_i(x)$$



## Kernel regression

$$\hat{y} = \frac{\sum_j w_j y_j}{\sum_j w_j}$$

$$w_j = K(x, x_j)$$
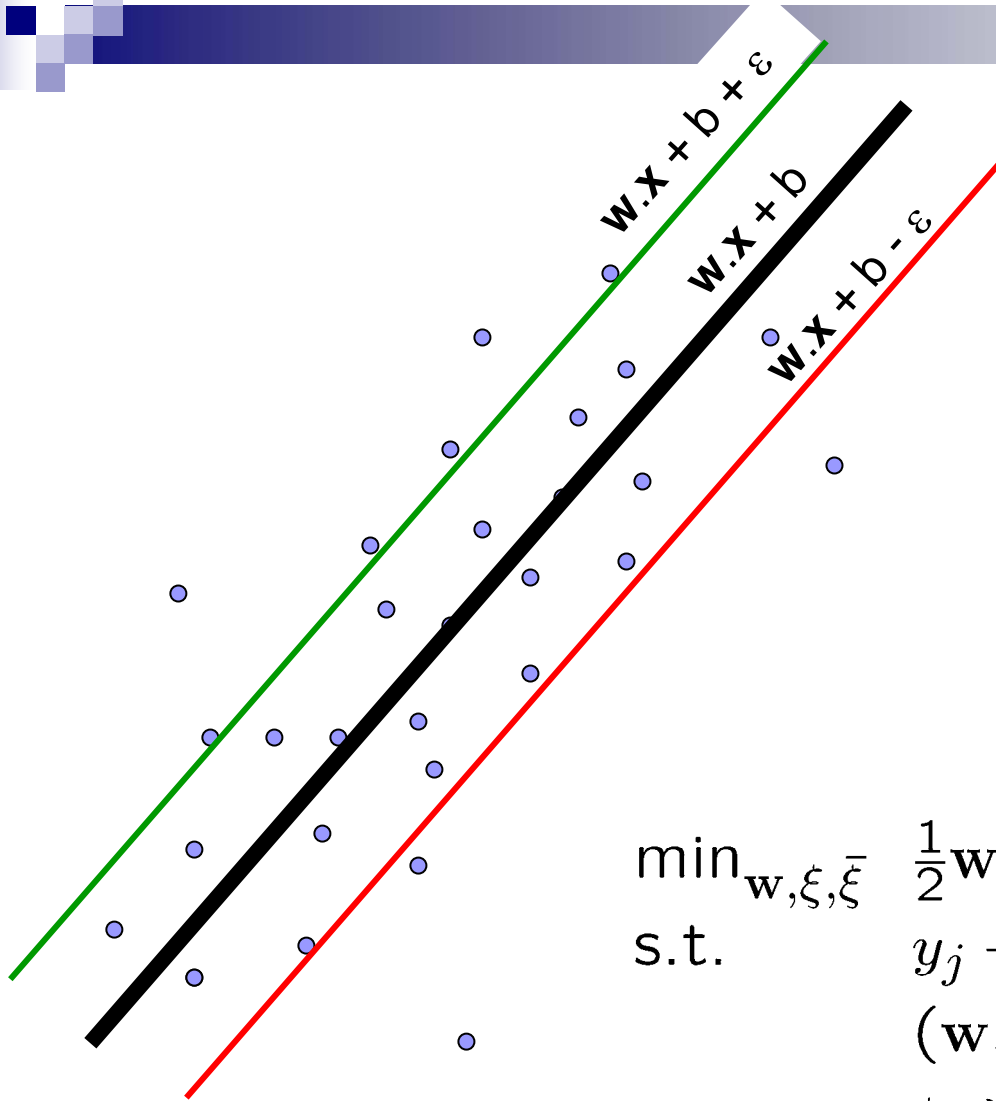
$K(\cdot, \cdot)$ has some parameters, gradient descent

## Kernel-weighted linear regression

Combine:
linear regression
with kernels

# Kernel-weighted linear regression

Local basis functions for each region



Kernels average between regions

# SVM regression



$\mathbf{w.x} + b + \epsilon$

$\mathbf{w.x} + b$

$\mathbf{w.x} + b - \epsilon$

$$\min_{\mathbf{w}, \xi, \bar{\xi}} \quad \frac{1}{2}\mathbf{w.w} + C \sum_{j=1}^{m}(\xi_j + \bar{\xi}_j)$$

$$\text{s.t.} \quad y_j - (\mathbf{w.x}_j + b) \leq \epsilon + \xi_j$$

$$(\mathbf{w.x}_j + b) - y_j \leq \epsilon + \bar{\xi}_j$$

$$\xi_j \geq 0, \quad \bar{\xi}_j \geq 0, \quad \forall j$$

# BIG PICTURE
## (a few points of comparison)

| | |
|---|---|
| DE | density estimation |
| CI | Classification |
| Reg | Regression |
| LL | Log-loss/MLE |
| Mrg | Margin-based |
| RMS | Squared error |

learning task

loss function

**Naïve Bayes**
DE, LL

**Boosting**
CI, exp-loss

iid Gaussian features

same H diff. loss

**Logistic regression**
DE, LL

**SVMs**
CI, Mrg

**SVM regression**
Reg, Mrg

**kernel regression**
Reg, RMS

**Instance-based Learning**
DE,CI,Reg

**Neural Nets**
DE,CI,Reg,RMS

**Decision trees**
DE,CI,Reg

**linear regression**
Reg, RMS

**This is a very incomplete view!!!**