

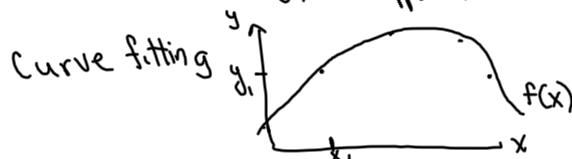
Dr. Cynthia Rudin (CCS, Columbia)

Regression (the Statistical Learning Version anyway):

Given $\{(x_i, y_i)\}_{i=1}^m$ where $x_i \in X$ & $y_i \in \mathbb{R}$, chosen randomly from an unknown probability distribution, find $f \in \mathcal{F}$

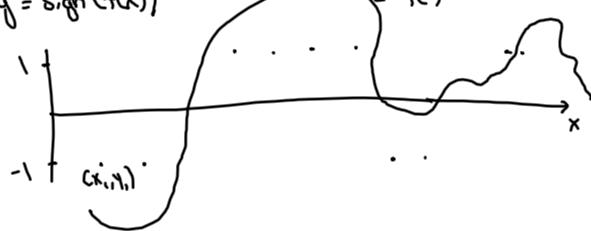
$f: X \rightarrow \mathbb{R}$, such that for a new randomly chosen instance $x \in X$, we have $f(x)$ close to its label y .

in some sense which varies on the application



Compare Classification & Regression

For classification, want to produce $f(x)$ such that

$$y = \text{sign}(f(x))$$


For regression, want to produce $f(x)$ such that

$$y = f(x)$$

Examples:

1) Want to estimate the number of customers for a given product, given data about the product and data about past product sales.

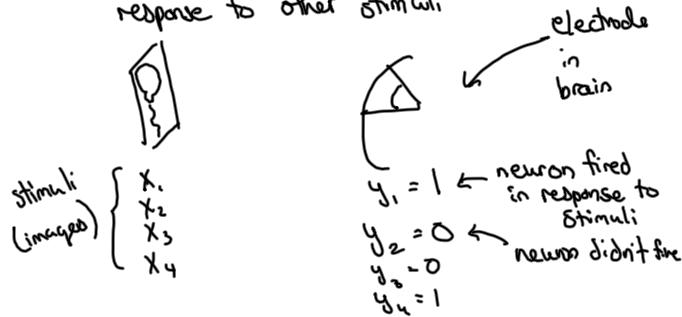
2) Time Series Analysis, for instance, predicting the price of stocks, or predicting the path of a moving object

3) Cancer Genomics (Hans & West 06)

Use gene expression levels (microarray data) from tumor samples to estimate:

- survival time
- amount of lymph node invasion
- response to a particular treatment

4) Computational Neuroscience: Calculating the firing rate of a neuron in the visual cortex for a new stimuli, based on the neuron's response to other stimuli



firing rate: prob that neuron will fire in response to stimulus x .

$$= P(y=1|x) = E(y|x) = \int y d\mu_{(x,y)}$$

(# of dimensions = # of pixels of image)

Square Loss and the Mean:

We generally would like to choose $f(x)$ so that it is close to y in some sense, for instance, we might want to choose f to minimize:

$$E[(y - f(x))^2 | x] \text{ for each } x$$

Say $y \in \{0, 1\}$. Fix x . What is the best $h = f(x)$? Let $p = P(y=1|x)$. Turns out the answer is $h = p$, the mean.

We cannot minimize $E[(y - f(x))^2]$ in practice, because we only have a finite sample of the data $\{(x_i, y_i)\}_{i=1..n}$

In practice, we minimize an "empirical" (based on data) quantity:

$$\sum_{i=1}^m (y_i - f(x_i))^2$$

<http://www.cs.columbia.edu/~allen/S08/NOTES/linkedlists.pdf>

Legendre published & named least squares in 1805 but Gauss claimed he had known the method since 1795. They both used it to determine from astronomical observation the orbits of bodies around the sun.

Solving Least Sq & Pseudo Inverse

Choose f to be linear. $f(\vec{x}) = \vec{\omega} \cdot \vec{x}$ for $\vec{x} \in \mathbb{R}^n$

Least Sq Problem: $\min_{\vec{\omega} \in \mathbb{R}^n} \sum_{i=1}^m (\vec{\omega} \cdot \vec{x}_i - y_i)^2$

Least J_{ϕ} : $\min_{\vec{\omega} \in \mathbb{R}^n} \sum_{i=1}^m (\vec{\omega} \cdot \vec{x}_i - y_i)^2$

$$\Phi(\vec{\omega})$$

To solve, set $\nabla \Phi(\vec{\omega}) = 0$

$$\nabla \Phi(\vec{\omega}) = \begin{pmatrix} \frac{d\Phi}{d\omega_1} \\ \vdots \\ \frac{d\Phi}{d\omega_n} \end{pmatrix} = \sum_{i=1}^m 2 \vec{x}_i^T (\vec{x}_i \cdot \vec{\omega} - y_i)$$

where $\vec{x} = \begin{pmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_m^T \end{pmatrix}$

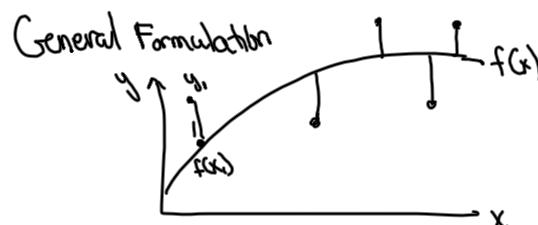
$$\nabla \Phi(\vec{\omega}) = 0 \Rightarrow \vec{\omega} = (\vec{x}^T \cdot \vec{x})^{-1} \vec{x}^T \vec{y}$$

pseudo inverse

If columns of \vec{x} are indep., pseudo inverse exists.

If \vec{x} has lin. indep. columns & rows, \vec{x}^{-1}
is the pseudoinverse

Pseudoinverse: Fredholm 1903
Moore 1920, Penrose 1955



For regression (& other learning problems) we often try to minimize a function of the form

$$L(f) = \sum_{i=1}^m l(y_i, f(x_i)) + c R(f)$$

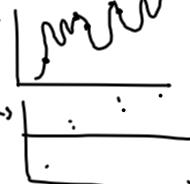
"empirical error" term
measures closeness of data to model

regularization parameter $c \in \mathbb{R}$

regularization term measured smoothness of f to prevent overfitting

If c is too small, model may overfit

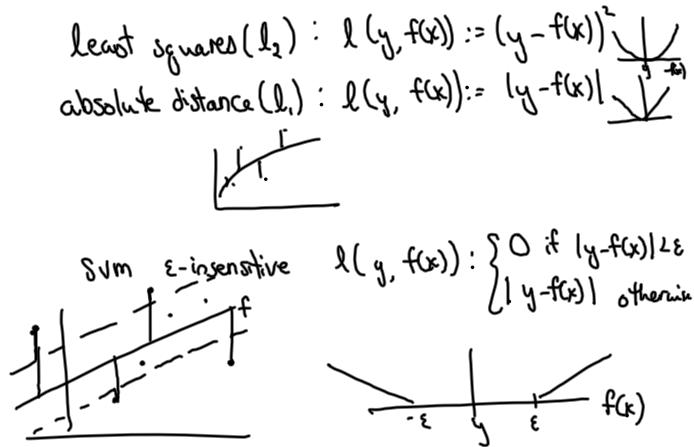
If c is too large, "underfit"



If c is just right, model generalizes (predicts) well



Some loss functions for regression:



Some regularization terms, assuming $f_w(\vec{x}) = \vec{w} \cdot \vec{x}$:

$$R(f_w) = \|\vec{w}\|_2^2 = \sum_j w_j^2$$

$$R(f_w) = \|\vec{w}\|_1 = \sum_j |w_j|$$

$$R(f_w) = \|f\|_{RKHS}^2 \quad (\text{you'll hear plenty about this later!})$$

Tikhonov Regularization Problem or Ridge Regression:

$$\min_{\vec{w}} \|\underline{\underline{X}} \vec{w} - \vec{y}\|_2^2 + \|\Gamma \vec{w}\|_2^2 \quad \text{can be solved:}$$

pos. semidef, $\underline{\underline{X}} \cdot \underline{\underline{X}}^T \geq 0$ forall

$$\vec{w} = (\underline{\underline{X}}^T \underline{\underline{X}} + \Gamma^T \Gamma)^{-1} \underline{\underline{X}}^T \vec{y}$$

True L0 sparsity

There has been a recent flurry of work (inspired by a result of Donoho) showing that in some cases (when columns of $\underline{\underline{X}}$ are almost orthogonal) the minimizer of:

$$\textcircled{1} \quad \|\underline{\underline{X}} \vec{w} - \vec{y}\|^2 + C \|\vec{w}\|_0$$

is the same as the minimizer of the real sparsity problem:

$$\textcircled{2} \quad \|\underline{\underline{X}} \vec{w} - \vec{y}\|^2 + C \|\vec{w}\|_0$$

where $\|\vec{w}\|_0$ = the number of nonzero elements of \vec{w} .

However $\textcircled{1}$ is convex & we can solve it whereas $\textcircled{2}$ is not!

Linear vs. Nonlinear Modelling form of the

There are many options for the function f .

$$f(\vec{x}) = \vec{w} \cdot \vec{x} \quad \text{linear}$$

$$f(\vec{x}) = \vec{w} \cdot \phi(\vec{x}) \quad \text{"non-linear"}$$

ϕ can be a vector in an infinite-dimensional Hilbert space.

This means w is infinite dimensional too!

This intro to regression is only the beginning!

Many other techniques:

- splines
- regression trees
- nnf networks & other nonlinear methods

In 1900 3 dims was a lot.

In 1980's 10 dims was a lot.

Now, we have ~10,000 dim