# Lecture 6. Regression

Consider learning the conditional distribution $p(y|x)$. This is often easier than learning the likelihoods $p(x|y)$, the prior $p(y)$, an then compute the posterior $p(y|x)$.
 Typically $x$ has many more dimensions than $y$.
 This is considered a regression problem. This is over 260 years old. Gauss finding the planetoid Ceres.

In this lecture, we adress three types of regression problem.

(I) Binary regression.    $y \in \{\pm 1\}$

$$p(y|\underline{x};\underline{\lambda}) = \frac{e^{y\underline{\lambda}\cdot\varphi(x)}}{e^{\underline{\lambda}\cdot\varphi(x)} + e^{-\underline{\lambda}\cdot\varphi(x)}} \ , \ \text{or} \ P(y|\underline{x};\underline{\lambda}) = \frac{e^{y\underline{\lambda}\cdot\varphi(x)}}{Z[\underline{x},\underline{\lambda}]}$$

$$\text{with } Z[\underline{x},\underline{\lambda}] = e^{\underline{\lambda}\cdot\varphi(x)} + e^{-\underline{\lambda}\cdot\varphi(x)}$$

(II)  Linear Gaussian

$$p(y|\underline{x};\underline{\lambda}) = \frac{1}{\sqrt{2\pi}\,\bar{\sigma}}\, e^{-\frac{(y - \underline{\lambda}\cdot\varphi(x))^2}{2\bar{\sigma}^2}} \quad (\text{linear in } \underline{\lambda})$$

This can be extended to vector-valued $y$.

In both cases, the parameters can be estimated by maximum likelihood (ML).

Minimize    $-\sum_{i=1}^{N} \log P(y_i|\underline{x}_i;\underline{\lambda})$
w.r.t $\underline{\lambda}$

This is a convex optimization problem.

We can also include a prior $P(\underline{\lambda})$ and perform MAP optimization.

(III)  Non-linear Regression:  E.G. Multi-Layer Perceptron

$$P(y|x) = \frac{1}{Z} e^{M(y, g(x;\underline{\lambda}))} \qquad (\text{Deep Networks})$$

$M(\cdot,\cdot)$ similarity measure
$g(x;\underline{\lambda})$ non-linear in $\underline{\lambda}$.

This leads to non-linear optimization problems
More powerful. More Computationally intensive.

(2)

## Binary Regression:

$$P(y|\underline{x}) = \frac{e^{y\,\underline{\lambda}\cdot\varphi(\underline{x})}}{e^{\underline{\lambda}\cdot\varphi(\underline{x})} + e^{-\underline{\lambda}\cdot\varphi(\underline{x})}}$$

Note: This leads to a decision rule:
classify $\underline{x}$ as $y = +1$, if $\underline{\lambda}\cdot\varphi(\underline{x}) > 0$
as $y = -1$, if $\underline{\lambda}\cdot\varphi(\underline{x}) < 0$
or $\hat{y}(\underline{x}) = \arg\min_y y\,\underline{\lambda}\cdot\varphi(\underline{x})$

Minimize:
$$-\sum_{i=1}^{N} \log P(y_i|\underline{x}_i; \underline{\lambda})$$

$$= -\sum_{i=1} y_i\,\underline{\lambda}\cdot\varphi(\underline{x}_i) + \sum_{i=1} \log\{ e^{\underline{\lambda}\cdot\varphi(\underline{x}_i)} + e^{-\underline{\lambda}\varphi(\underline{x}_i)} \}$$

This is a convex function of $\underline{\lambda}$
( check → second term can be written as $\sum \underline{\lambda}\cdot\varphi(\underline{x}_i)$ )
$+ \sum \log\{ 1 + e^{-\underline{\lambda}\cdot\varphi(\underline{x})} \}$... Hessian is +ve definite

The gradient of $-\sum_{i=1}^{N} \log P(y_i|\underline{x}_i; \underline{\lambda})$ w.r.t. $\underline{\lambda}$
is $-\sum_{i=1}^{N} y_i\,\varphi(\underline{x}_i) + \sum_{i=1} \sum_{y\in\{\pm1\}} y\,\varphi(\underline{x}_i) P(y|\underline{x}_i; \underline{\lambda})$

Hence the minimum also balances the statistics: $\hat{\underline{\lambda}}$ s.t.

$$\sum_{i=1}^{N} y_i\,\varphi(\underline{x}_i) = \sum_{i=1}^{N} \sum_{y\in\{\pm1\}} y\,\varphi(\underline{x}_i) P(y|\underline{x}_i; \underline{\lambda}).$$

Steepest Descent: $\underline{\lambda}^{t+1} = \underline{\lambda}^t - \Delta\{ -\sum_{i=1}^{N} y_i\,\phi(\underline{x}_i) + \sum_{i=1} \sum_{y\in\{\pm1\}} y\,\phi(\underline{x}_i) P(y|\underline{x}_i, \underline{\lambda}) \}$
or variants, can be used
to minimize $-\sum_{i=1} \log P(y_i|\underline{x}_i; \underline{\lambda})$

Special case — idealized "neuron"

$$\varphi(\underline{x}) = (x_1, x_2 \ldots x_n)$$
$$\underline{\lambda} = (\lambda_1, \lambda_2, \ldots \lambda_n)$$

$$\underline{\lambda}\cdot\varphi(\underline{x}) = \lambda_1 x_1 + \lambda_2 x_2 + \ldots + \lambda_n x_n.$$

Integrate - And - Fire: If $\underline{\lambda}\cdot\varphi(\underline{x}) > 0$, then the
neuron will fire with probability.
or threshold $\underline{\lambda} = (\lambda_0, \lambda_1, \ldots \lambda_n)$
$\varphi(\underline{x}) = (1, x_1, \ldots x_n)$
$$P(y|x) = \frac{e^{y\,\underline{\lambda}\cdot\varphi(\underline{x})}}{e^{\underline{\lambda}\cdot\phi(\underline{x})} + e^{-\underline{\lambda}\phi(\underline{x})}}$$
$\underline{\lambda}\cdot\varphi(\underline{x}) > 0$, $\sum_{i=1} \lambda_i x_i > -\lambda_0$.

# (II) Gaussian Linear Regression.

$$y = \underline{\lambda} \cdot \underline{\varphi}(x) + \epsilon \qquad P(\epsilon) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\epsilon^2}{2\sigma^2}}$$

zero mean Gaussian

$$P(y \mid \underline{x}, \underline{\lambda}, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y - \underline{\lambda} \cdot \underline{\varphi}(\underline{x}))^2}$$

**ML:** minimize $-\sum_{i=1}^{N} \log P(y_i \mid \underline{x}_i, \underline{\lambda}, \sigma)$

dataset $\chi = \{(\underline{x}_i, y_i) : i = 1, .., N\}$

$$= \sum_{i=1}^{N} \frac{1}{2\sigma^2} (y_i - \underline{\lambda} \cdot \underline{\varphi}(\underline{x}_i))^2 - N \log \{\sqrt{2\pi}\sigma\}.$$

<u>Minimize</u> w.r.t. $\underline{\lambda}$

$$-\frac{1}{\sigma^2} \sum_{i=1}^{N} (y_i - \underline{\lambda} \cdot \underline{\varphi}(\underline{x}_i)) \, \underline{\varphi}(\underline{x}_i) = 0$$

Solution $\quad \underline{\hat{\lambda}} = \left\{ \sum_{i=1}^{N} \underline{\varphi}(\underline{x}_i) \, \underline{\varphi}(\underline{x}_i)^T \right\}^{-1} \sum_{i=1}^{N} y_i \, \underline{\varphi}(\underline{x}_i)$

<u>Minimize</u> w.r.t. $\sigma$ $\quad -\frac{1}{\sigma^3} \sum_{i=1}^{N} (y_i - \underline{\lambda} \cdot \underline{\varphi}(\underline{x}_i))^2 - \frac{N}{\sigma}$

Solution $\quad \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^{N} (y_i - \underline{\hat{\lambda}} \cdot \underline{\varphi}(\underline{x}_i))^2$ //

This estimates the regression coefficients $\underline{\lambda}$ and also the variance.

This can be generalized to allow for vector-valued output.

<u>Note:</u> in terms of coordinates

$$\sum_{i=1}^{N} (y_i - \sum_a \lambda_a \varphi_a(\underline{x}_i))(y_i - \sum_b \lambda_b \varphi_b(\underline{x}_i))$$

$$\frac{\partial}{\partial \lambda_a} = \quad 2(-1) \sum_{i=1}^{N} \varphi_a(\underline{x}_i)(y_i - \sum_b \lambda_b \varphi_b(\underline{x}_i))$$

Then the minimum occurs at

$$\sum_b \left\{ \sum_{i=1}^{N} \varphi_a(\underline{x}_i) \varphi_b(\underline{x}_i) \right\} \lambda_b = \sum_{i=1}^{N} y_i \varphi_a(\underline{x}_i).$$

$\sum_{i=1}^{N} \varphi_a(\underline{x}_i) \varphi_b(\underline{x}_i)$ are the coefficients of the $a^{th}$ row and $b^{th}$ column of a matrix.

## $L^1 \cdot$ Variant

$$y = \underline{\lambda} \cdot \varphi(\underline{x}) + \epsilon \quad , \quad \text{where} \quad P(\epsilon) = \frac{1}{2\sigma} e^{-|\epsilon|/\sigma}$$

$$\left( \int_0^\infty e^{-\epsilon/\sigma} d\epsilon = \left[ -\sigma e^{-\frac{\epsilon}{\sigma}} \right]_0^\infty = \sigma \right).$$

$$P(y | \underline{x}, \underline{\lambda}, \sigma) = \frac{1}{2\sigma} e^{-\frac{1}{\sigma}|y - \underline{\lambda} \cdot \varphi(\underline{x})|}.$$

Estimate $\underline{\lambda}, \sigma$ by ML

$$\frac{1}{\sigma} \sum_{i=1}^{N} |y_i - \underline{\lambda} \cdot \varphi(\underline{x}_i)| + N \log(2\sigma)$$

$$\hat{\underline{\lambda}} = \underset{\underline{\lambda}}{\text{ARG MIN}} \sum_{i=1}^{N} |y_i - \underline{\lambda} \cdot \varphi(\underline{x}_i)|$$

This is a convex optimization problem — steepest descent and other algorithms will get $\hat{\underline{\lambda}}$

$$\hat{\sigma} = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{\underline{\lambda}} \cdot \varphi(\underline{x}_i)|.$$

This is more robust than the previous (Gaussian) method. It is $L_1$ norm instead of $L_2$.

(5) ## Specific Examples:

$$\underline{\lambda} \cdot \underline{\phi}(\underline{x}) = \omega_0 + \omega_1 x_1 + \ldots + \omega_d x_d = \sum_{j=1}^{d} \omega_j x_j + \omega_0$$

$$\underline{\lambda} = (\omega_0, \omega_1, \ldots, \omega_d), \quad \underline{\phi}(\underline{x}) = (1, x_1, \ldots, x_d)$$

Simplest Case: $\qquad y = \omega_1 x + \omega_0$

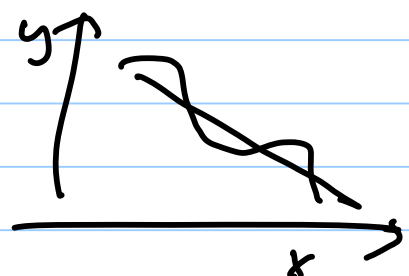Dataset $\chi_\nu = \langle (x^i, y^i) : i = 1 \text{ to } N \rangle$

$$E(\omega_1, \omega_0 \mid \chi_N) = \sum_{i=1}^{N} \left\{ y^i - (\omega_1 x^i + \omega_0) \right\}^2$$

Minimize : $\dfrac{\partial E}{\partial \omega_1} = 0$ & $\dfrac{\partial E}{\partial \omega_0} = 0$

solution $\begin{cases} \widehat{\omega}_1 = \dfrac{\sum\limits_{i=1}^{N} x^i y^i - \bar{x}\bar{y} N}{\sum\limits_i (x^i)^2 - N\bar{x}^2} \\[20pt] \widehat{\omega}_0 = \bar{y} - \omega_1 \bar{x} \end{cases}$

where $\bar{x} = \sum\limits_{i=1}^{N} x_i / N, \quad \bar{y} = \sum\limits_i y_i / N$



A "richer" model can be used.

$\rightarrow e.g \quad \underline{\lambda} \cdot \underline{\phi}(\underline{x}) = \omega_2 x^2 + \omega_1 x + \omega_0$

Too high an order follows the data too closely.

(6)          More abstractly

Linear regression: $\quad\quad\quad\quad 1. \phi(\underline{x}) = \omega_1 x^i + \omega_0$.

Differentiate energy wrt $\omega_0, \omega_0$ gives two equations

$$\sum_i y^i = N\omega_0 + \omega_1 \sum_i x^i$$

$$\sum_i g_i x^i = \omega_0 \sum_i x^i + \omega_1 \sum_i (x^i)^2.$$

Expressed in linear algebra form as $\underline{A}\,\underline{\omega} = \underline{y}$

$$\underline{A} = \begin{bmatrix} N & \sum_i x^i \\ \sum_i y^i & \sum_i (x^i)^2 \end{bmatrix}, \quad \underline{\omega} = \begin{bmatrix} \omega_0 \\ \omega_1 \end{bmatrix}, \quad \underline{y} = \begin{bmatrix} \sum_i y^i \\ \sum_i x^i y^i \end{bmatrix}$$

Solved to give $\quad \underline{\omega} = \underline{A}^{-1} \underline{y}$.

Polynomial Regression.

$$g(x^i | \omega_k \ldots, \omega_2, \omega_1, \omega_0) = \omega_k (x^i)^k + \ldots + \omega_1 x^i + \omega_0.$$

$k+1$ parameters $\omega_k \ldots \omega_0$

Diff. energy — gives $k+1$ linear eqn's in $k+1$ variables.

$$\underline{A}\,\underline{\omega} = \underline{y}.$$

Can write $\underline{A} = \underline{D}^T \underline{D}$, $\quad \underline{y} = \underline{D}^T \underline{r}$ $\quad \underline{D} = \begin{bmatrix} 1 & x_1 & x_1^k \\ 1 & x_2 & x_2^k \\ 1 & \cdots & \cdots \end{bmatrix}$

Solve to get $\underline{\omega} = (D^T D)^{-1} D^T r$ $\quad \underline{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \end{bmatrix}$

Must adjust the complexity of the model to the amount of data available
    Complexity of poly regression is no. parameters $k$.
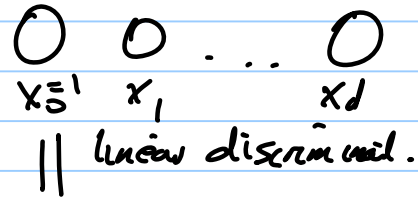Need to pick $k$ to give best generalization error

# Multilayer Perceptrons

Analogy to Neural Networks in the Brain.

$\qquad$ — over-simplified.

Perceptron.
$$y = \sum_{j=1}^{d} \omega_j x_j + \omega_0$$

idealized neuron

$x_0 = 1 \quad x_1 \quad \ldots \quad x_d$

$\|$ linear discriminant.

threshold function.

hard
$$S(a) = \begin{cases} 1, & \text{if } a > 0 \\ 0, & \text{otherwise.} \end{cases}$$

soft.
$$y = \bar{\sigma}(\underline{\omega}^T \underline{x}) = \frac{1}{1 + e^{-\underline{\omega}^T \underline{x}}}.$$

$\bar{\sigma}(.)$ sigmoid function.

There are a variety of different algorithms to train a perceptron from labelled examples

Example: Quadratic error.
$$E(\underline{\omega} \mid \underline{x}^t, y^t) = \frac{1}{2}(y^t - \underline{\omega} \cdot \underline{x}^t)^2$$

update rule
$$\Delta \omega_j^t = -\Delta \frac{\partial E}{\partial \omega_j} = +\Delta(y^t - \underline{\omega} \underline{x}^t) \underline{x}^t$$

$$E(\{\underline{\omega}_i\} \mid \underline{x}^t, y^s) = -\sum_i \{r_i^t \log y_i^t + (1 - r_i^t) \log(1 - y_i^t)\}$$

$$r^t = \text{sigmoid}(\omega^T{}_i \underline{x}^t).$$

update rule
$$\Delta \omega_j^t = -\eta(r^t - y^t) x_j^t.$$

Update = Learning Factor · (Desired Output − Actual Output) × Input.

# Multilayer Perceptron.

A single layer perceptron can only approximate linear function of the input — i.e. the set of perceptrons has limited capacity.

Multilayer perceptrons were invented to increase the capacity — introduce hidden units (nodes)

$$z_h = \sigma(\underline{\omega}_h^T \underline{x})$$

$$= \frac{1}{1 + \exp\left\{-\sum_{j=1}^{d} \omega_{hj} x_j + \omega_{i0}\right\}},$$

$h = 1, \ldots H.$

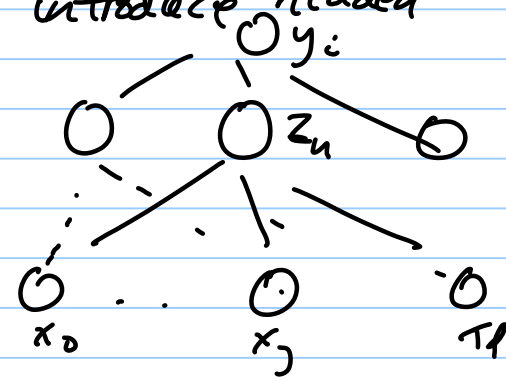Output $y_i = \underline{v}_i^T z = \sum_{h=1}^{H} v_{ih} z_h + v_{i0}.$

Other output function — e.g. $y_i = \sigma(\underline{v}_i^T \underline{z}),$

## Many levels can be specified.

What do the hidden units represent?
Unclear, but many people have tried to explain them.
Any input-output function can be represented as a multilayer perceptron with enough hidden units
~ infinite capacity.

# How to train a multilayer perceptron?

Unknown parameters — the weights $\omega_{hj}, v_{ij}$.

Define an error function:

eg. $\quad E[\omega, v] = \sum_i \left( y_i - \sum_h v_{ih} \; \sigma \; (\sum_j \omega_{hj} x_j) \right)^2$

**update** $\quad \Delta \omega_{hj} = -\dfrac{\partial E}{\partial \omega_{hj}}$  <span style="color:red">computed by the chain rule.</span>

$\quad\quad\quad\quad \Delta v_{ik} = \dfrac{-\partial E}{\partial v_{ih}}$  <span style="color:red">computed directly</span>

<span style="color:red">Define $\quad r_k = \sigma\left(\sum_j \omega_{kj} x_j\right), \quad E = \sum_i \left(y_i - \sum_k v_{ik} r_k\right)^2$</span>

$\dfrac{\partial E}{\partial \omega_{kj}} = \dfrac{\partial E}{\partial r_k} \cdot \dfrac{\partial r_k}{\partial \omega_{kj}}$

$\dfrac{\partial E}{\partial r_k} = -2 \sum_i \left(y_i - \sum_\ell v_{i\ell} r_\ell\right) v_{ik}$

$\dfrac{\partial r_k}{\partial \omega_{kj}} = x_j \, \sigma'\left(\sum_j \omega_{kj} x_j\right)$

$\sigma'(z) = \dfrac{d}{dz} \sigma(z) = \sigma(z)\langle 1 - \sigma(z)\rangle.$

<span style="color:red">Hence</span> $\quad \dfrac{\partial E}{\partial \omega_{hj}} = -2 \sum_i \left(y_i - \sum_\ell v_{i\ell} r_\ell\right) v_{ik} \, x_k \, \sigma\left(\sum_j \omega_{kj} x_j\right)$
$\langle 1 - \sigma\left(\sum_j \omega_{hj} x_j\right)\rangle$

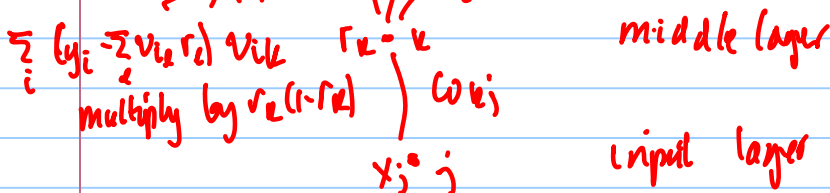<span style="color:red">error at output layer</span>     <span style="color:red">weights from middle layer to output layer.</span>

This is called **backpropagation**. The error at the output layer is propagated back to the nodes at the middle layer

$\sum_i \left(y_i - \sum_\ell v_{i\ell} r_\ell\right) v_{ih}$

<span style="color:red">where it is multiplied by the activity $r_k(1-r_k)$ at that node, and by the activity $x_j$ at the input.</span>

<span style="color:red">→ errors</span>

<span style="color:red">$\sum_i (y_i - \sum_\ell v_{i\ell} r_\ell) v_{ik}$   $r_{k-k}$</span>

<span style="color:red">multiply by $r_k(1-r_k)$ $\Big)$ $\omega_{kj}$</span>

<span style="color:red">$x_j \cdot j$</span>

<span style="color:red">middle layer</span>

<span style="color:red">input layer</span>

## Learning in Batch Mode:

Put all data into an energy function - i.e. Sum the errors over all the training data. Update the weights — equations above — by summing over all the data.

Alternatively, online learning.

At time $t$, select one element $(\underline{x}^t, y^t)$ from the training set at random.
Perform one iteration of steepest descent only. Then select another element at random.

Online learning can be shown to converge and may be better at avoiding local minima in the energy function than batch methods.
Also, online is suitable if we keep getting new input over time — this happens in many real world applications.

Critical Issues for multilayer perceptrons
— how many hidden units to use?

Book describes several techniques for dealing with this → for example, having more hidden units than you need and then penalizing the weights.
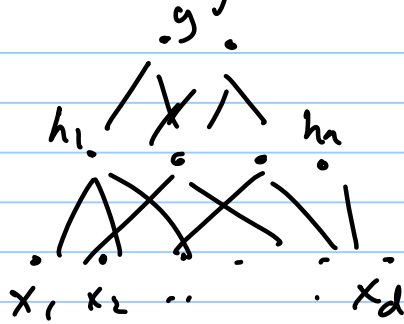
E.G. Add term $\sum_{h,j} (\omega_{hj}^2) + \sum_{i,h} (v_{ih}^2)$ to the cost function → intuition, if the weights are small to a hidden unit, then the hidden unit is not used.

In practice, some of the most effective multi layer perceptions are those which the structure was hand designed.

# Multilayer Perceptrons / SVM / AdaBoost

$$y_i = \sum_h v_{ih} h_h$$

$$h_h = \sum_j \bar{\sigma}\left(\sum \omega_{hj} x_j\right).$$

<u>SVM</u> can also be represented in this way.

$$y = \text{sign}\left(\sum_\mu \alpha_\mu y_\mu \, \underline{X}_\mu \cdot \underline{X}\right)$$

hidden units response  $\underline{X}_\mu \cdot \underline{X} = h_\mu.$

$$y = \text{sign}\left(\sum_\mu \alpha_\mu y_\mu h_\mu\right)$$

Advantage of SVM — number of hidden units is given by the no. of support vectors.

→ $\langle \alpha_\mu \rangle$ specified by minimizing the primal problem (well defined algorithm to perform this minimization).