

(1.)

Lecture 6 . Regression

Spring 2013.

Note Title

4/21/2013

Consider learning the conditional distribution $p(y|x)$.
This is often easier than learning the likelihoods $p(x|y)$, the prior $p(y)$, and 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 address three types of regression problem.

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

$$p(y|x;\lambda) = \frac{e^{y\lambda \cdot \phi(x)}}{e^{\lambda \cdot \phi(x)} + e^{-\lambda \cdot \phi(x)}}, \text{ or } P(y|x;\lambda) = \frac{e^{y\lambda \cdot \phi(x)}}{z(x;\lambda)}$$

with $z(x;\lambda) = e^{\lambda \cdot \phi(x)} + e^{-\lambda \cdot \phi(x)}$

(II) Gaussian regression - $(y - \lambda \cdot \phi(x))^2$

$$p(y|x;\lambda) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y - \lambda \cdot \phi(x))^2}{2\sigma^2}}$$

This can be extended to vector-valued y .

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

Minimize $-\sum_{i=1}^n \log P(y_i|x_i;\lambda)$
wrt λ

This is a convex optimization problem.

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

(III) Multi-Layer Perception.

This is an extension of (I).

It leads to non-convex optimization.

Binary Regression:

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

Note: this leads to a decision rule:

classify x as $y = +1$, if $\lambda \cdot \phi(x) > 0$
as $y = -1$, if $\lambda \cdot \phi(x) < 0$
or $\hat{y}(x) = \arg \min_y y \lambda \cdot \phi(x)$

(2)

Spring 2013.

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

$$= -\sum_{i=1}^N y_i \lambda \cdot \varphi(x_i) + \sum_{i=1}^N \log \left(e^{\lambda \cdot \varphi(x_i)} + e^{-\lambda \cdot \varphi(x_i)} \right)$$

This is a convex function of λ
 (check \rightarrow second term can be written as $\sum \lambda \cdot \varphi(x_i)$)
 $+ \sum \log \{ 1 + e^{-2\lambda \cdot \varphi(x_i)} \}$... Hessian is the definite

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

Hence the minimum also balances the s-statistics: $\hat{\lambda}$ st.

$$\sum_{i=1}^N y_i \varphi(x_i) = \sum_{i=1}^N \sum_{y \in \{+1, -1\}} y \varphi(x_i) P(y | x_i; \hat{\lambda})$$

Note: also can use a prior $p(\lambda)$. (later lecture).

Algorithms - steepest descent, or variational -

can be specified to minimize this

Special case - idealized "neuron"

$$\varphi(\underline{x}) = (x_1, x_2, \dots, x_n)$$

$$\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)$$

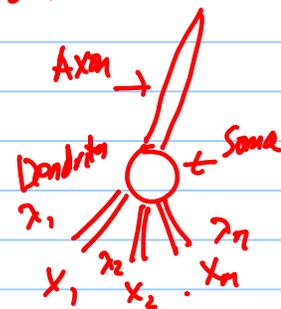
$$\lambda \cdot \varphi(\underline{x}) = \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n$$

Integrate-and-Fire: if $\lambda \cdot \varphi(\underline{x}) > 0$, then the neuron will fire with probability.

or threshold $\lambda = (\lambda_0, \lambda_1, \dots, \lambda_n)$

$$\varphi(\underline{x}) = (1, x_1, \dots, x_n)$$

$$\lambda \cdot \varphi(\underline{x}) > 0 \quad , \quad \sum_{i=1}^n \lambda_i x_i > -\lambda_0$$



(3)

Spring 2013.

Now Consider continuous regression with Gaussian model.

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

3rd mean Gaussian

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

ML: minimize $-\sum_{i=1}^N \log P(y_i | x_i, \underline{\lambda}, \sigma)$

$$= \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - \underline{\lambda} \cdot \underline{\phi}(x_i))^2 - N \log \{ \sqrt{2\pi}\sigma \}$$

dataset $\mathcal{X} = \{(x_i, y_i) : i=1, \dots, N\}$

Minimize w.r.t. $\underline{\lambda}$

$$-\frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \underline{\lambda} \cdot \underline{\phi}(x_i)) \underline{\phi}(x_i) = 0$$

Solution $\hat{\underline{\lambda}} = \left(\sum_{i=1}^N \underline{\phi}(x_i) \underline{\phi}(x_i)^T \right)^{-1} \sum_{i=1}^N y_i \underline{\phi}(x_i)$

Minimize w.r.t. σ

$$-\frac{1}{\sigma^3} \sum_{i=1}^N (y_i - \underline{\lambda} \cdot \underline{\phi}(x_i))^2 - \frac{N}{\sigma}$$

Solution $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\underline{\lambda}} \cdot \underline{\phi}(x_i))^2$ //

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

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

Similarly binary regression can be extended to multiclass: replace $y \underline{\lambda} \cdot \underline{\phi}(x)$ by $\underline{\lambda} \cdot \underline{\phi}(y, x)$.

Note: in terms of coordinates

$$\frac{\partial}{\partial \lambda_a} = \sum_{i=1}^N (y_i - \sum_b \lambda_b \phi_b(x_i)) \phi_a(x_i)$$

Then the minimum occurs at

$$\sum_b \left\{ \sum_{i=1}^N \phi_a(x_i) \phi_b(x_i) \right\} \lambda_b = \sum_{i=1}^N y_i \phi_a(x_i)$$

$\sum_{i=1}^N \phi_a(x_i) \phi_b(x_i)$ are the coefficients of the a^{th} row and b^{th} column of a matrix.

(4)

Spring 2013.

We can get a variant by specifying:

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

$$\left(\int_0^\infty e^{-\epsilon/\sigma} d\epsilon = \left[-\sigma e^{-\epsilon/\sigma} \right]_0^\infty = \sigma \right)^{\frac{1}{\sqrt{2\sigma}}}$$

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

Estimate λ, σ by ML

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

$$\hat{\lambda} = \underset{\lambda}{\text{Arg Min}} \sum_{i=1}^N |y_i - \lambda \cdot \varphi(x_i)|$$

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

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

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

Next we consider some specific examples, from Alpaydm's book.

Spring 2013

(5)

Alypudin's notation is different.

He uses a function $g(\underline{x}) = \underline{\lambda} \cdot \underline{\varphi}(\underline{x})$
(in our notation).

For example,

$$\underline{\lambda} = (\omega_0, \omega_1, \dots, \omega_d)$$

$$g(\underline{x}) = \omega_1 x_1 + \dots + \omega_d x_d + \omega_0 = \sum_{j=1}^d \omega_j x_j + \omega_0$$

Back to simple case - single attribute x ,

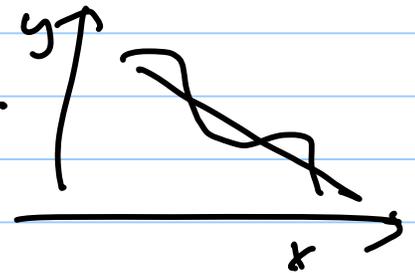
Dataset $X = \{(x^i, y^i) : i=1, \dots, N\}$

$$E(\omega_1, \omega_0 | X) = \sum_{i=1}^N (y^i - (\omega_1 x^i + \omega_0))^2$$

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

solution
$$\begin{cases} \hat{\omega}_1 = \frac{\sum_{i=1}^N x^i y^i - \bar{x} \bar{y} N}{\sum_{i=1}^N (x^i)^2 - N \bar{x}^2} \\ \hat{\omega}_0 = \bar{y} - \hat{\omega}_1 \bar{x} \end{cases}$$

where $\bar{x} = \frac{\sum_{i=1}^N x^i}{N}$, $\bar{y} = \frac{\sum_{i=1}^N y^i}{N}$



A "richer" model can be used.

→ e.g. $g(x) = \omega_2 x^2 + \omega_1 x + \omega_0$

Too high an order follows the data too closely.

(6)

Spring 2013

More Abstractly,

Linear regression: $g(x^i | \omega_1, \omega_0) = \omega_1 x^i + \omega_0$.

Differentiate energy w.r.t ω_1, ω_0 gives two equations

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

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

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

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

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

Polynomial Regression.

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

$k+1$ parameters $\omega_k, \dots, \omega_0$

Diff. energy - gives $k+1$ linear eq's in $k+1$ variables.

$$\underline{A}\underline{w} = \underline{y}$$

Can write $\underline{A} = \underline{D}^T \underline{D}$, $\underline{y} = \underline{D}^T \underline{r}$

Solve to get $\underline{w} = (\underline{D}^T \underline{D})^{-1} \underline{D}^T \underline{r}$

$$\underline{D} = \begin{bmatrix} x_1 & x_1^k \\ x_2 & x_2^k \\ \vdots & \vdots \end{bmatrix}, \quad \underline{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \end{bmatrix}$$

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

(7)

Multilayer Perceptrons

Spring 2013.

Note Title

4/27/2008

Analogy to Neural Networks in the Brain.
- over-simplified.

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

idealized neuron

hard threshold function.

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

linear discriminant.

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

$\sigma(\cdot)$ sigmoid function.

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

Example: Quadratic error.

$$E(\underline{\omega} | \underline{x}^t, y^t) = \frac{1}{2} (y^t - \underline{\omega} \cdot \underline{x}^t)^2$$

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

$$E(\{\omega_i\} | \underline{x}^t, y^t) = -\sum_i (r_i^t \log y_i^t + (1-r_i^t) \log(1-y_i^t))$$

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

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

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

Multilayer Perceptron.

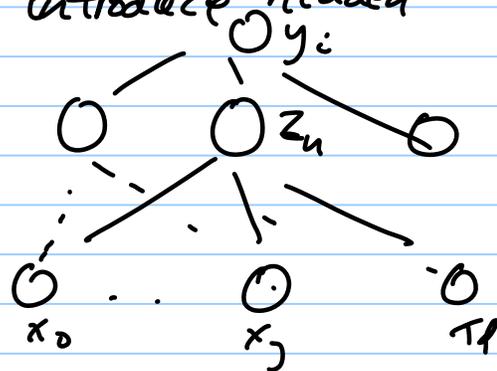
A single layer perceptron can only approximate linear functions of the input — i.e. the set of perceptrons has limited capacity.

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

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

$$= \frac{1}{1 + \exp\left(-\sum_{j=1}^d \omega_{hj} x_j + \omega_{h0}\right)}$$

$h=1, \dots, H.$



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

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

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.

(9)

Spring 13

How to train a multilayer perceptron?

Unknown parameters - the weights w_{hj}, v_{ij} .

Define an error function:

e.g.
$$E[w, v] = \sum_i (y_i - \sum_h v_{ih} \sigma(\sum_j w_{hj} x_j))^2$$

update
$$\Delta w_{hj} = -\frac{\partial E}{\partial w_{hj}}$$
 Computed by the chain rule.

$$\Delta v_{ik} = -\frac{\partial E}{\partial v_{ik}}$$
 Computed directly

Define $r_k = \sigma(\sum_j w_{kj} x_j)$, $E = \sum_i (y_i - \sum_k v_{ik} r_k)^2$

$$\frac{\partial E}{\partial w_{kj}} = \frac{\partial E}{\partial r_k} \cdot \frac{\partial r_k}{\partial w_{kj}}$$

$$\frac{\partial E}{\partial r_k} = -2 \sum_i (y_i - \sum_l v_{il} r_l) v_{ik}$$

$$\frac{\partial r_k}{\partial w_{kj}} = x_j \sigma'(\sum_j w_{kj} x_j)$$

$$\sigma'(z) = \frac{d}{dz} \sigma(z) = \sigma(z) (1 - \sigma(z))$$

Here
$$\frac{\partial E}{\partial w_{hj}} = -2 \sum_i (y_i - \sum_l v_{il} r_l) v_{ih} x_j \sigma(\sum_j w_{hj} x_j) (1 - \sigma(\sum_j w_{hj} x_j))$$

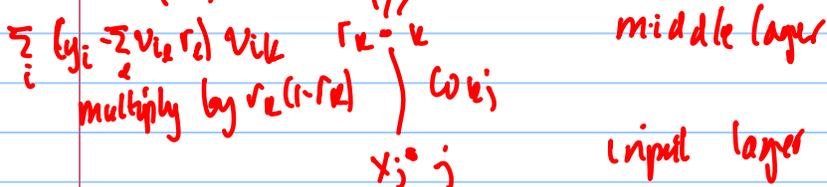
error at output layer weights from middle layer to output layer.

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

$$\sum_i (y_i - \sum_l v_{il} r_l) v_{ih}$$

where it is multiplied by the activity $r_k (1 - r_k)$ at that node, and by the activity x_j at the input.

→ errors δ_k



(10)

Spring 2013

Learning in Batch Mode:

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

Alternatively, online learning.

At time t , select one element (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.

(11)

Spring 2013

Critical Issues for multilayer perceptrons

— how many hidden units to use?

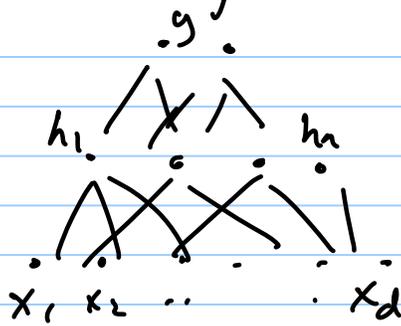
Book describes several techniques for dealing with this → for example, having more hidden units than you need and then penalizing the weights.

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 perceptrons are those which the structure was hand designed.

(12)

Multilayer Perceptrons / SVM / AdaBoost Spring 2017 (next two lectures)



$$y_i = \sum_h v_{ih} h_h$$
$$h_h = \sum_j \sigma(\sum_j \omega_{hj} x_j)$$

SVM 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.

→ $\{d_n\}$ specified by minimizing the primal problem (well defined algorithm to perform this minimization).