# **Machine Learning**

# **Linear Regression**

Lesson 3

## **Linear Regression**

- Basics of Regression
- Least Squares estimation
- Polynomial Regression
- Basis functions Regression model
- Regularized Regression
- Statistical Regression
  - Maximum Likelihood (ML) estimation
  - Maximum A-Posteriori (MAP) estimation
  - Bayesian Regression

## **Basics in Regression**

- Input: Set of N tuples:  $D = \{(x_1, t_1), \dots, (x_N, t_N)\}$ -  $x_i \in \mathbb{R}^d$  is the sample
  - − t<sub>i</sub> ∈ R prediction or a value of an unknown function over the features of x<sub>i</sub>
- Supervised technique
- **Goal:** create a function  $\mathbf{y}: \forall x_i \in D : y(x_i, \theta) \approx t_i$ 
  - in one dimension  $y:\mathbb{R}\to\mathbb{R}$
  - in d dimensions  $y: \mathbb{R}^D \to \mathbb{R}$

#### **\theta** is the (*unknown*) set of parameters

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – **ML1**(3)

### **Graphical Example of Regression**







## Linear regression: The case of 1-dimensional data



## Linear regression: The case of 1-dimensional data



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (8)

### Linear regression: The case of 1-dimensional data



Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (9)

### Least squares linear fit to data

- Most popular estimation method is least squares.
- Determine linear weights w that minimize the sum of squared loss (SSL):

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} (y(x_i, w) - t_i)^2 = \sum_{i=1}^{N} (w_0 + w_1 x_i - t_i)^2$$

- Use standard differential calculus:
  - differentiate SSL with respect to  $w_0$ ,  $w_1$
  - find zeros of each partial differential equation
  - solve for  $w_0$ ,  $w_1$

• Derivatives of parameters:

$$\frac{\partial J(w)}{\partial w_0} = \sum_{i=1}^N w_0 + w_1 x_i - t_i = 0 \Longrightarrow N w_0 + w_1 \sum_{i=1}^N x_i = \sum_{i=1}^N t_i$$
$$\frac{\partial J(w)}{\partial w_1} = \sum_{i=1}^N (w_0 + w_1 x_i - t_i) x_i = 0 \Longrightarrow w_0 \sum_{i=1}^N x_i + w_1 \sum_{i=1}^N x_i^2 = \sum_{i=1}^N x_i t_i$$

$$\hat{w}_{1} = \frac{\frac{1}{N} \sum_{i=1}^{N} x_{i} t_{i} - (\bar{x})(\bar{t})}{\frac{1}{x^{2}} - (\bar{x})^{2}} = \frac{\operatorname{cov}(x, t)}{\operatorname{var}(x)} \quad \hat{w}_{0} = \bar{t} - \hat{w}_{1} \bar{x}$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (11)

#### **More dimensions (d>1)** $x \in R^d$





- Input x has *d* **features**:  $\vec{x} = \begin{bmatrix} x_1 \dots x_d \end{bmatrix}$
- There are d+1 linear weights for describing the regression function y(x, w)

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (12)

Linear regression model

$$y(x,w) = w_0 + w_1 x_1 + \ldots + w_d x_d = w_0 + \sum_{j=1}^d w_j x_j$$

Alternative representation

$$y(x,w) = w^T \widetilde{x}$$

where 
$$W = \begin{bmatrix} w_0 & w_1 & \dots & w_d \end{bmatrix}$$
  
 $\widetilde{x} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$ 

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (13)

## Sum of Squared Error (SSE)

How can we quantify the error?

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} (y(x_i, w) - t_i)^2$$

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} \left( w_0 + w^T x_i - t_i \right)^2 = \frac{1}{2} \sum_{i=1}^{N} \left( w^T \widetilde{x}_i - t_i \right)^2$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (14)

### Sum of Squared Error (SSE)

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} (y(x_i, w) - t_i)^2 = \frac{1}{2} \sum_{i=1}^{N} (w^T \tilde{x}_i - t_i)^2$$
  
**Dbservation**  

$$y(x, w)$$
  
**Prediction**  

$$y(x, w) - t = (w^T \tilde{x} - t)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (15)

• How can we quantify the error?

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \hat{x}_{i} - t_{i} \right)^{2} = \frac{1}{2} \left\| Xw - T \right\|^{2}$$

$$J(w) = \frac{1}{2} (Xw - T)^T (Xw - T)$$



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (16)

## Finding good parameters

- Want to find parameters w which minimize the error
- Think of a cost "surface": error residual for that **w**...



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (17)

# **SSE Minimization**

- Consider a simple problem
  - One feature (d=1), two data points (N>2)
  - Two unknowns:  $W_0$ ,  $W_1$
  - Two equations:

 $t_1 = w_0 + w_1 x_1$ 

$$t_2 = w_0 + w_1 x_2$$



• Can solve this system directly (X: 2x2):

$$T = Xw \implies \hat{w} = X^{-1}T$$

- However, most of the time, N > d+1
  - There may be no linear function that hits all the data exactly
  - Instead, solve directly for minimum of SSE function

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (18)

# The general case



$$J(w) = \frac{1}{2} \|Xw - T\|^{2} = \frac{1}{2} (Xw - T)^{T} (Xw - T)$$

$$\nabla_{w}J(w) = 0 \Longrightarrow X^{T}(Xw - T) = 0$$

• Reordering, we take

$$X^{T}X\hat{w} = X^{T}T \Longrightarrow \hat{w}_{LS} = (X^{T}X)^{-1}X^{T}T$$
  
Least Squares estimator

- (X<sup>T</sup> X)<sup>-1</sup>X<sup>T</sup> is called the "pseudo-inverse"
- If  $X^{T}$  is square and independent, this is the inverse  $X^{-1}$

#### **Optimization methods of SSE**

- Even when (X<sup>T</sup> X) is invertible, might be computationally expensive if X is huge.
- Treat as an **optimization problem**:

$$\hat{w} = \arg\min_{w} J(w) = \arg\min_{w} \frac{1}{2} \|Xw - T\|^2$$

• How to find an estimator?



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (20)

#### **Gradient Descent Optimization**

#### <u>Steepest Descent</u>

Function decreases most quickly in the direction of the negative gradient.



*n*: learning rate (or step-size)

$$w^{(new)} = w^{(old)} - \eta \frac{\partial J(w)}{\partial w}$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (21)

## Effect of step-size (n)



- Large η : fast convergence but larger residual error. May cause oscillations
- **Small η :** slow convergence but small residual error.

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (22)

#### **1. Gradient descent Optimization scheme**

• Since *J(w)* is convex, move along negative of gradient



Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (23)

## 2. Stochastic Gradient descent

Previous scheme was a batch gradient descent: all training examples are participated at every step

$$w^{(new)} = w^{(old)} - \eta X^{T} (Xw^{(old)} - T) = w^{(old)} - \eta \sum_{i=1}^{N} (w^{(old)^{T}} x_{i} - t_{i}) x_{i}$$

 Stochastic gradient descent scheme repeatedly examines a single example at every step:

for 
$$i = 1, ..., N$$
  

$$w^{(new)} = w^{(old)} - \eta \left( w^{(old)^T} x_i - t_i \right) x_i$$

- Advantage: often gets w close to the minimum much faster than batch mode.
- It is preferred over batch when training set is large.

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (24)

#### 3. Newton-Raphson Optimization scheme

- Update rule:  $w^{(new)} = w^{(old)} H^{-1}\nabla J(w)$ where H: Hessian matrix (second derivatives of J(w))  $\nabla J(w) = X^T X w - X^T T$  $H = \nabla \nabla J(w) = X^T X$
- By substituting the previous we obtain the rule:

$$w^{(new)} = w^{(old)} - (X^T X)^{-1} (X^T X w^{(old)} - X^T T) = (X^T X)^{-1} X^T T$$

• This is the least-squares solution and is exact solution in one step

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (25)

#### Effects of SSE (12 error) choice

Sensitivity to outliers



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (26)

#### Use the sum of absolute error (SAE) (11 error)



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (27)

#### Use the sum of absolute error (SAE) (11 error)



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (28)

#### Use the sum of absolute error (SAE) (11 error)



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (29)

## **Non-linear Regression**

- Consider non-linear regression
  - Ex: higher-order polynomials





<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (30)

#### **Polynomial Regression**

$$y(x,w) = w_0 + w_1 x + \dots + w_\rho x^\rho = w_0 + \sum_{k=1}^{\rho} w_k x^k = w^T \phi(x)$$

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2} = \frac{1}{2} \left\| \Phi w - T \right\|^{2}$$

$$\Phi = \begin{bmatrix} 1 & x_1 & \cdots & x_1^{\rho} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \cdots & x_N^{\rho} \end{bmatrix} \qquad T = \begin{bmatrix} t_1 & \cdots & t_N \end{bmatrix}^T$$
$$w = \begin{bmatrix} w_0 & w_1 & \cdots & w_d \end{bmatrix}^T$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (31)

• Example:

Use 3d polynomial. Single feature x, predict target t:

Sometimes useful to think of "feature transform"

$$y(x) = w^T \phi(x)$$
  $\phi(x) = [1, x, x^2, x^3]$ 

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (32)

#### Learning polynomial regression coefficients

Using Least-squares:

$$\hat{w}_{LS} = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{T}$$

• Using Gradient descent update rule:

$$w^{(new)} = w^{(old)} - \eta \Phi^T \left( \Phi w^{(old)} - T \right)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (33)

## **Linear Basis Function Models**

• In general, can use **any feature set** we think is useful

$$\phi(x) = [\phi_1(x), \dots, \phi_m(x)]$$

- Other information about the problem
  - e.g. location, age, ...
- Polynomial functions
  - Features [1, x, x<sup>2</sup>, x<sup>3</sup>, ...]
- Other functions
  - 1/x, sqrt(x),  $x_1 * x_2$ , ...

$$y(x) = \sum_{j=1}^{m} w_j \phi_j(x) = w^T \phi(x)$$

#### • Regression remains "Linear" = linear in the parameters

Features we can make as complex as we want!

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (34)

 $\phi_j(x): R^d \to R$ basis functions

## **Examples of Basis Function**

#### **Polynomial basis functions**

$$\phi_j(x) = x^j.$$

- These are **global**
- A small change in *x* affect all basis functions.



#### **Gaussian basis functions**

$$\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$$

- These are **local**
- A small change in X only affect nearby basis functions.
- Parameters  $\mu_j$  and **s** control location and scale (width).
- Related to kernel methods.


### **Sigmoid basis functions**

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

- Also these are local: a small change in x only affect nearby basis functions.
- Parameters  $\mu_j$  and  $\mathbf{s}$  control location and scale (slope).



# Additive models

- The basis functions can capture various properties of the inputs (e.g. qualitative)
- These are called "Additive models"
- For example: we can try to rate documents based on text descriptors

$$\mathbf{x} = \text{text document (collection of words)}$$
  

$$\phi_i(\mathbf{x}) = \begin{cases} 1 \text{ if word } i | \text{ appears in the document} \\ 0 \text{ otherwise} \end{cases}$$
  

$$f(\mathbf{x}; \mathbf{w}) = w_0 + \sum_{i \in \text{words}} w_i \phi_i(\mathbf{x})$$

 We can view the additive models graphically in terms of simple "units" and "weights".



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (39)

# Overfitting



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (40)

# Overfitting



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – **ML1** (41)

# **Overfitting and complexity**

- More complex models will always fit the training data better
- But they may "overfit" the training data, learning complex relationships that are not really present



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (42)

# Training versus test error

- Plot MSE as a function of model complexity
  - Polynomial order
- Decreases
  - More complex function fits training data better
- What about new data?
- Low order
  - Error decreases
  - Underfitting
- Higher order
  - Error increases
  - Overfitting



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – **ML1** (43)

# **Dealing with Overfitting**

- ✓ Use more data
- ✓Use a tuning set
- ✓ Regularization
- ✓ Be a Bayesian

44

# **1. Avoiding overfitting: Cross-validation**

 Cross-validation allows us to estimate the generalization error based on training examples alone



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (45)

• Leave-one-out cross-validation treats each training example in turn as a test example:

$$CV = \frac{1}{N} \sum_{i=1}^{N} \left( y(x_i, \hat{w}^{-i}) - t_i \right)^2$$

where  $\hat{w}^{-i}$  are the least squares estimates of the parameters without the *i*<sup>th</sup> training example.

### **Polynomial Regression example**



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (47)

# 2. Regularization

- Overfitting is reduced by imposing a constraint on the overall magnitude of the parameters.
- Objective function is modified

$$J(w) = E_D(w) + \lambda E_W(w)$$

**Data term + Regularization term** 

- $\lambda$  : regularization parameter
- **E**<sub>w</sub>(w) set constraints to linear weights

# **L2-Regularization or Ridge Regression**

• The regularization term is quadratic

$$E_{W}(w) = \frac{1}{2} w^{T} w = \frac{1}{2} \sum_{j=1}^{M} w_{j}^{2}$$

that penalize large weights

Objective function

$$J(w) = \frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2} + \frac{\lambda}{2} w^{T} w =$$
$$= \frac{1}{2} \left\| \Phi w - T \right\|^{2} + \frac{\lambda}{2} w^{T} w$$

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (49)

$$\nabla_{W} J(w) = 0 \Longrightarrow \Phi^{\mathsf{T}} \Phi \hat{w} - \Phi^{\mathsf{T}} T + \lambda \hat{w} = 0 \Longrightarrow$$

$$\hat{w} = (\Phi^{T} \Phi + \lambda I) \Phi^{T} T$$
Identity matrix (m x m)

- Ridge regression or weight decay
- Since the squared weights is compatible with the squared error function, we get a nice closed form solution for the optimal weights.

# **More general regularizers**

• Several regularized regression models



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (51)

# L1 Regularization (LASSO)

• The regularization term is

$$E_{W}(w) = \frac{1}{2} \|w\|_{1} = \frac{1}{2} \sum_{j=1}^{M} |w_{j}|$$

Objective function

$$J(w) = \frac{1}{2} \|\Phi w - T\|_{2}^{2} + \frac{\lambda}{2} \|w\|_{1}$$

 Ability to create sparse models that are more easily interpreted

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (52)

### **Constrained formulation of the L1 Regularization**

$$\min_{w} \left\{ \frac{1}{2} \| \Phi w - T \|_{2}^{2} \right\} \quad s.t. \| w \|_{1} = 0$$

- Least Absolute Selection and Shrinkage Operator (LASSO)
- Several optimization techniques for solving the problem

# An Optimization Scheme: Sequentially added sign constraints

$$\min_{w} \left\{ \frac{1}{2} \| Xw - T \|_{2}^{2} \right\} \qquad s.t. \| w \|_{1} \le L \qquad \qquad \begin{array}{c} & \stackrel{+w_{1} + w_{2} + w_{3} \le t}{+w_{1} - w_{2} - w_{3} \le t} \\ & \stackrel{+w_{1} - w_{2} - w_{3} \le t}{+w_{1} - w_{2} - w_{3} \le t} \\ & \stackrel{-w_{1} + w_{2} + w_{3} \le t}{-w_{1} + w_{2} + w_{3} \le t} \end{array}$$

 Examine all possible combinations of the sign -w1-w2-w3 ≤ 1 -w1-w2-w3 ≤ 1
 of the elements of w.

#### Algorithm 1 Tibshirani's Method

1: 
$$w = LS(X, y)$$

- 2:  $constraints = \{null\}$
- 3: while  $||w||_1 \le L \, do$
- 4: add sign(w) to constraints

**Tibshirani, R.** (1996). Regression shrinkage and selection via the lasso. Journal of Royal. Statist. Soc B., Vol. 58, No. 1, pages 267-288).

 $|a_{1}|$   $|a_{2}|$   $|a_{2}|$   $|a_{2}|$   $|a_{2}|$  t

5:  $w = min_w ||Xw - y||_2^2$  subject to constraints

#### 6: end while

### L2 vs. L1 regularization



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (55)

### L2 vs. L1 regularization

.

L2 regularized regression	L1 regularized regression
Quadratic regularization has the advantage that the <b>solution is in closed form</b> (not appeared in non-quadratic regularizers).	Lasso represents a <b>convex</b> <b>optimization problem</b> solved by quadratic programming or other convex optimization methods
L2 regularization <b>shrinks</b> <b>coefficients towards (but not</b> <b>to) zero</b> , and towards each other.	L1 regularization shrinks coefficients to zero at different rates; different values of $\lambda$ give models with different subsets of features.

# **Model Likelihood and Empirical Risk**

Two related but distinct ways to look at a model.

– Empirical Risk: "How much error does the model have on the training data?"

– Model Likelihood: "What is the likelihood that a model generated the observed data?"

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (57)

# Statistical view of linear regression

• Assumption: (Generative model)

**Observed output = constructive function + stochastic noise** 

$$t = y(x, w) + \varepsilon$$

 Whatever we cannot capture with our chosen family of functions will be interpreted as noise



# Statistical view of linear regression

$$t = y(x, w) + \varepsilon$$

• If we consider (white) Gaussian noise

$$\mathcal{E} \sim N(0, \beta^{-1})$$
 where  $\beta = \frac{1}{\sigma^2}$  precision or inverse variance

then we introduce the stochastic model:

$$t \mid x \sim N(y(x, w), \beta^{-1})$$

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (59)

# **Model Likelihood**

- Training set:  $D = \{(x_1, t_1), \dots, (x_N, t_N)\}$
- Generative model:  $\theta = \{w, \beta\}$

$$p(t_i \mid x_i, \theta) = N(t_i \mid y(x_i, w), \beta^{-1}) = \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left\{-\frac{\beta}{2}(y(x_i, w) - t_i)^2\right\}$$

• Likelihood function:

$$p(D | \theta) = \prod_{i=1}^{N} p(t_i | x_i, \theta) = \prod_{i=1}^{N} N(t_i | w^T \phi(x_i), \beta^{-1})$$

assuming Independently Identically distributed (*iid*) data

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (60)

• Log-likelihood:

$$L(\theta) = \ln p(D \mid \theta) = \sum_{i=1}^{N} \ln p(t_i \mid x_i, \theta)$$

$$L(\theta) = -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{i=1}^{N} (y(x_i, w) - t_i)^2$$

• Assuming a linear regression model of m basis functions

$$L(\theta = \{w, \beta\}) = -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{i=1}^{N} (w^{T} \phi(x_{i}) - t_{i})^{2}$$

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (61)

Maximum Likelihood (ML) estimation

$$L(\theta = \{w, \theta\}) = -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{i=1}^{N} (w^{T} \phi(x_{i}) - t_{i})^{2}$$

Partial derivative of linear weights (w):

$$\frac{\partial L(\theta)}{\partial w} = -\frac{\beta}{2} \sum_{i=1}^{N} \left( \hat{w}^{T} \phi(x_{i}) - t_{i} \right) \phi^{T}(x_{i}) = 0$$

$$\hat{w}^T \sum_{i=1}^N \phi(x_i) \phi^T(x_i) - \sum_{i=1}^N t_i \phi^T(x_i) = 0 \Longrightarrow \Phi^T \Phi \hat{w} = \Phi^T T$$

$$\hat{w}_{ML} = \left(\Phi^T \Phi\right)^{-1} \Phi^T T = \hat{w}_{LS}$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (62)

Maximum Likelihood (ML) estimation

$$L(\theta = \{w, \theta\}) = -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{i=1}^{N} (w^{T} \phi(x_{i}) - t_{i})^{2}$$

• Partial derivative of inverse variance (**B**):

$$\frac{\partial L(\theta)}{\partial \beta} = -\frac{N}{2\hat{\beta}} - \frac{1}{2} \sum_{i=1}^{N} \left( \widehat{w}^T \phi(x_i) - t_i \right)^2 = 0$$

$$\frac{1}{\hat{\beta}_{ML}} = \frac{1}{N} \sum_{i=1}^{N} \left( \widehat{w}_{ML}^{T} \phi(x_{i}) - t_{i} \right)^{2} = \frac{1}{N} \left\| \Phi \widehat{w}_{ML} - T \right\|_{2}^{2}$$

#### Mean prediction squared error

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (63)

# **MAP** estimation of Linear Regression

 MAP derivation of linear regression assumes a prior distribution over linear weights w:

$$p(w \mid \alpha) = N(w \mid 0, \alpha^{-1}I) = \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\left\{-\frac{\alpha}{2}w^Tw\right\}$$

**α**: is a **hyperparameter** over **w** and is the *precision* or *inverse variance* of the distribution.

• Then, the **posterior distribution** is obtained as

$$p(w|X,T,\alpha,\beta) \propto p(T|X,w,\beta)p(w|\alpha)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (64)

Optimize the Bayesian posterior  $p(w | X, T, \alpha, \beta) \propto p(T | X, w, \beta) p(w | \alpha)$ 

• Set the MAP log-likelihood function:

$$L_{MAP}(\theta) = \ln p(T \mid X, w, \beta) + \ln p(w \mid \alpha)$$

likelihood + prior

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (65)

Optimize the Bayesian posterior  $p(w \mid X, T, \alpha, \beta) \propto p(T \mid X, w, \beta) p(w \mid \alpha)$ 

• Set the MAP log-likelihood function:

$$L_{MAP}(\theta) = \ln p(T \mid X, w, \beta) + \ln p(w \mid \alpha)$$

likelihood

$$L_{ML}(\theta) = \ln p(T \mid X, w, \beta) = \sum_{i=1}^{N} \ln p(t_i \mid x_i, w, \beta) =$$
$$= -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{i=1}^{N} (w^T \phi(x_i) - t_i)^2$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (66)

Optimize the Bayesian posterior  $p(w | X, T, \alpha, \beta) \propto p(T | X, w, \beta) p(w | \alpha)$ 

• Set the MAP log-likelihood function:

$$L_{MAP}(\theta) = \ln p(T \mid X, w, \beta) + \ln p(w \mid \alpha)$$

prior

$$\ln p(w \mid \alpha) = \ln N(w \mid 0, \alpha^{-1}I) =$$
$$= -\frac{M}{2}\ln 2\pi + \frac{M}{2}\ln \alpha - \frac{\alpha}{2}w^{T}w$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (67)

#### Optimize the Bayesian posterior

• Ignoring terms that do not depend on w :

$$L_{MAP}(w) = \ln p(T \mid X, w, \beta) + \ln p(w \mid \alpha) =$$
  
$$\approx -\frac{\beta}{2} \sum_{i=1}^{N} (w^T \phi(x_i) - t_i)^2 - \frac{\alpha}{2} w^T w$$

• MAP estimation of linear weights:

$$\frac{\partial L_{MAP}(w)}{\partial w} = 0 \Longrightarrow \hat{w}_{MAP} = \left(\Phi^T \Phi + \lambda I\right) \Phi^T T$$

"Thus regularized (ridge) L2 regression reflects a 0-mean isotropic Gaussian prior on the weights"

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (68)

### A summary of the Linear Regression techniques

#### **Deterministic** approaches

Least squares

$$\frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2}$$

**Statistical approaches** 

Maximum Likelihood (ML)  $\varepsilon \sim N(0, \beta^{-1})$  $p(t | w, \beta) = N(w^T \phi(x), \beta^{-1})$ 

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (69)

### A summary of the Linear Regression techniques

Deterministic	approaches

Least squares

$$\frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2}$$

#### L2 regularized (ridge)

$$\frac{1}{2}\sum_{i=1}^{N} \left(w^{T}\phi(x_{i}) - t_{i}\right)^{2} + \frac{\lambda}{2}\sum_{j=1}^{m} \left(w_{j}\right)^{2}$$

**Statistical approaches** 

Maximum Likelihood (ML)  $\varepsilon \sim N(0, \beta^{-1})$  $p(t \mid w, \beta) = N(w^T \phi(x), \beta^{-1})$ 

**MAP with Gaussian prior**  $p(w | \alpha) = N(w | 0, \alpha^{-1}I)$ 

$$\ln p(T \mid X, w, \beta) + \ln p(w \mid \alpha)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (70)

### A summary of the Linear Regression techniques

Deterministic	c approaches

Least squares

$$\frac{1}{2} \sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2}$$

#### L2 regularized (ridge)

$$\frac{1}{2}\sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2} + \frac{\lambda}{2} \sum_{j=1}^{m} \left( w_{j} \right)^{2}$$

L1 regularized (lasso)

$$\frac{1}{2}\sum_{i=1}^{N} \left( w^{T} \phi(x_{i}) - t_{i} \right)^{2} + \frac{\lambda}{2} \sum_{j=1}^{m} \left| w_{j} \right|$$

#### **Statistical approaches**

Maximum Likelihood (ML)  $\varepsilon \sim N(0, \beta^{-1})$  $p(t \mid w, \beta) = N(w^T \phi(x), \beta^{-1})$ 

**MAP with Gaussian prior**  $p(w | \alpha) = N(w | 0, \alpha^{-1}I)$ 

$$\ln p(T | X, w, \beta) + \ln p(w | \alpha)$$

**MAP with Laplacian prior** 

 $w \sim Ce^{-\lambda |w|_1}$ 



<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (71)

# **Hierarchical Bayesian model**



Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (72)
## **Alternative Regression models**

Elastic Net – Complex Prior for the weights w

$$w \sim C e^{-\lambda_1 |w|_1 + \lambda_2 |w|_2}$$

Penalization by weighted L1 and L2 norms

$$J(w) = \frac{1}{2} \left\| \Phi w - T \right\|_{2}^{2} + \frac{\lambda_{1}}{2} \left\| w \right\|_{1} + \frac{\lambda_{2}}{2} \left\| w \right\|_{2}^{2}$$

➤ Weighted Least Squares: Assign for each case (x<sub>i</sub>,t<sub>i</sub>) a weight v<sub>i</sub> ≥ 0 (the higher the more important the case)

$$\frac{1}{2} \sum_{i=1}^{N} v_i \left( w^T \phi(x_i) - t_i \right)^2 = \left\| V^{\frac{1}{2}} (\Phi w - T) \right\|^2$$
  
Then:  $\hat{w}_{WLS} = \left( \Phi^T V \Phi \right)^{-1} \Phi^T V T$   $V = \begin{bmatrix} v_1 & v_2 & 0 \\ v_2 & 0 \\ 0 & v_N \end{bmatrix}$ 

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (73)

## **Bayesian Linear Regression**

- The previous MLE or MAP derivation of linear regression uses **point estimates** for the weight vector, w.
- Bayesian modeling estimates the **posterior** distribution of weights after receiving all observations.
- This allows us to find the distribution of the target of the new coming input, and thus make **prediction**.

- Generative model:  $t = w^T \phi(x) + \varepsilon$   $p(t \mid x, w, \beta) = N(t \mid w^T \phi(x), \beta^{-1})$ • Let N observations:  $D = \{(x_1, t_1), \dots, (x_N, t_N)\}$
- Then, the join distribution of N observations  $T = \{t_1, ..., t_N\}$  is:

$$p(T | X, w) = \prod_{i=1}^{N} p(t_i | x_i, w) = N(T | \Phi w, \beta^{-1}I)$$

We treat that linear weights w as Gaussian random variables with mean m<sub>0</sub> and covariance matrix S<sub>0</sub>

$$p(w) = N(m_0, S_0)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (75)

• Posterior distribution: product of two Gaussians

$$p(w|T,X) \propto p(T|X,w)p(w) = N(T|\Phi w,\beta^{-1}I)N(w|m_0,S_0)$$

If we have a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 (B.42)

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}+\mathbf{b},\mathbf{L}^{-1})$$
 (B.43)

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (B.44)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\},\mathbf{\Sigma})$$
 (B.45)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}. \tag{B.46}$$

C. M. Bishop, "Pattern Recognition and Machine Learning", page 689

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (76)

 Substituting we take the general form of the posterior distribution:

$$p(w|T,X) = N(T | \Phi w, \beta^{-1}I)N(w|m_0, S_0) = N(w|m_N, S_N)$$

where:

$$S_N = (\beta \Phi^T \Phi + S_0^{-1})^{-1}$$
  $m_N = S_N (\beta \Phi^T T + S_0^{-1} m_0)$ 

• If  $m_0 = 0$  and  $S_0 = \alpha^{-1}$  / then

$$S_N = \left(\beta \Phi^T \Phi + aI\right)^{-1} \qquad \lambda = \frac{\alpha}{\beta}$$

$$m_N = \left(\beta \Phi^T \Phi + \alpha I\right)^{-1} \beta \Phi^T T = \left(\Phi^T \Phi + \lambda I\right)^{-1} \Phi^T T$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (77)

## **Predictive distribution**

 Assuming a new input x<sup>\*</sup> we are looking for making a prediction of its target t<sup>\*</sup>. This is equivalent on estimating the posterior distribution:

$$p(t^* \mid x^*, T)$$

• Marginalizing we take:

$$p(t^* | x^*, T) = \int p(t^* | x^*, w) p(w | T) dw$$

where  

$$p(w|T) = N(w|m_N, S_N)$$
 $m_N = (\beta \Phi^T \Phi + \alpha I)^{-1} \beta \Phi^T T$ 
 $S_N = (\beta \Phi^T \Phi + \alpha I)^{-1}$ 
 $p(t^* | x^*, w) = N(t^* | w^T \phi(x^*), \beta^{-1})$ 

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML1 (78)

**Predictive distribution (cont.)** 

$$p(t^* | x^*, T) = \int p(t^* | x^*, w) p(w | T) dw$$
$$p(w | T) = N(w | m_N, S_N) \qquad S_N = (\beta \Phi^T \Phi + aI)^{-1}$$
$$p(t^* | x^*, w) = N(t^* | w^T \phi(x^*), \beta^{-1})$$

• According to Gaussian properties (**B.44**) we receive:

$$p(t^* | x^*, T) = N(t^* | m_N^T \phi(x^*), \sigma_N^2(x^*))$$

where:

$$\sigma_N^2(x^*) = \beta^{-1} + \phi^T(x^*) S_N \phi(x^*)$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML1 (79)