Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 4 - Linear Regression - Probabilistic Interpretation and
Regularization

- Need to determine **w** for the linear function $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^{n} w_i \phi_i(\mathbf{x_j}) = \mathbf{\Phi w}$ which minimizes our error function $E(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$
- Owing to basis function $\phi$, "Linear Regression" is *linear* in **w** but NOT in **x** (which could be arbitrarily non-linear)!

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x_1}) & \phi_2(\mathbf{x_1}) & ...... & \phi_p(\mathbf{x_1}) \\ \cdot & & & \\ \cdot & & & \\ \phi_1(\mathbf{x_m}) & \phi_2(\mathbf{x_m}) & ...... & \phi_n(\mathbf{x_m}) \end{bmatrix} \quad (1)$$

*function* $f(x, w)$     *in* $x$

- Need to determine $\mathbf{w}$ for the linear function $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^{n} w_i \phi_i(\mathbf{x_j}) = \mathbf{\Phi w}$ which minimizes our error function $E(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$
- Owing to basis function $\phi$, "Linear Regression" is *linear* in $\mathbf{w}$ but NOT in $\mathbf{x}$ (which could be arbitrarily non-linear)!

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x_1}) & \phi_2(\mathbf{x_1}) & ...... & \phi_p(\mathbf{x_1}) \\ \cdot & & & \\ \cdot & & & \\ \phi_1(\mathbf{x_m}) & \phi_2(\mathbf{x_m}) & ...... & \phi_n(\mathbf{x_m}) \end{bmatrix} \tag{1}$$

- Least Squares error and corresponding estimates:

$$E^* = \min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \min_{\mathbf{w}} \left( \mathbf{w^T \Phi^T \Phi w} - 2\mathbf{y^T \Phi w} + \mathbf{y^T y} \right) \tag{2}$$
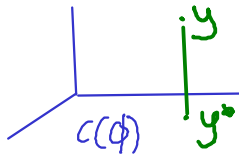
*Derived graphically*

$$\mathbf{w^*} = \arg\min_{\mathbf{w}} \mathbf{E}(\mathbf{w}, \mathcal{D}) = \arg\min_{\mathbf{w}} \left\{ \sum_{j=1}^{m} \left( \sum_{i=1}^{n} \mathbf{w_i} \phi_i(\mathbf{x_j}) - \mathbf{y_j} \right)^2 \right\} \tag{3}$$

- Let $\mathbf{y}^*$ be a solution in the column space of $\Phi$
- The least squares solution is such that the distance between $\mathbf{y}^*$ and $\mathbf{y}$ is minimized
- Therefore, the line joining $\mathbf{y}^*$ to $\mathbf{y}$ should be orthogonal to the column space of $\Phi$
  $\Rightarrow$

$$\mathbf{w} = (\mathbf{\Phi^T\Phi})^{-1}\mathbf{\Phi^Ty} \tag{4}$$

- Here $\Phi^T\Phi$ is invertible only if $\Phi$ has full column rank

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori Estimates, Regularization
3. How to minimize the resultant and more complex error functions?
   - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

- Linear Model: $Y$ is a linear function of $\phi(x)$, subject to a random noise variable $\varepsilon$ which we believe is 'mostly' bounded by some threshold $\sigma$:

$$Y = w^T\phi(x) + \varepsilon$$
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

(exponential, $\chi^2$, t, uniform)

- Motivation: $\mathcal{N}(\mu, \sigma^2)$, has maximum entropy among all real-valued distributions with a specified variance $\sigma^2$

- $3 - \sigma$ rule: About 68% of values drawn from $\mathcal{N}(\mu, \sigma^2)$ are within one standard deviation $\sigma$ away from the mean $\mu$; about 95% of the values lie within $2\sigma$; and about 99.7% are within $3\sigma$.

$\overset{*}{p}$ is a distr to be estimated

$P$ is family of pdfs $\supseteq \{N,$ expnential, $t, x^2 \dots\}$

$$\overset{*}{p} = \underset{}{\text{maximize scope}} \equiv \underset{p}{\max} \int_x \left(-\log_2 p(x)\right) p(x) \, dx$$

$$\text{s.t } \mathrm{var}_p[X] = \sigma^{-2}$$

<span style="color:magenta">Budgeted encoding using $P$</span>

<span style="color:red">more prob $\Rightarrow$ less specified</span>

distr over outcomes $\{$day, hour, mins, secs $\dots\}$

<span style="color:red">More prob $\Rightarrow$ more imp</span>   <span style="color:red">increasing imp</span>

Q: What will $p^*$ be if $\mathcal{P}$ is family of discrete distributions?

$$\underset{p_1, p_2 \ldots p_k}{\arg\max} \sum_i -p_i \log_2 p_i = \left\{ p_1^* = p_2^* \cdots = p_k^* = \frac{1}{k} \right\}$$

Why isn't uniform the entropy maximizer in the continuous case?

Ans:



pdf=0
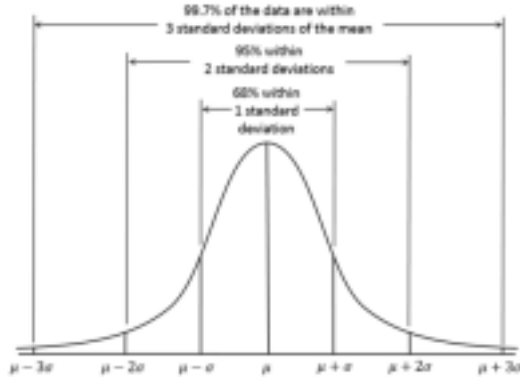
$\leftarrow$ pdf $\frac{1}{x_2 - x_1}$

pdf=0

$x_1$          $x_2$

Figure 1: $3 - \sigma$ rule: About 68% of values drawn from $\mathcal{N}(\mu, \sigma^2)$ are within one standard deviation $\sigma$ away from the mean $\mu$; about 95% of the values lie within $2\sigma$; and about 99.7% are within $3\sigma$. Source: https://en.wikipedia.org/wiki/Normal_distribution

- Linear Model: $Y$ is a linear function of $\phi(\mathbf{x})$, subject to a random noise variable $\varepsilon$ which we believe is 'mostly' around some threshold $\sigma$:

$$Y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$$
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- This allows for the Probabilistic model

$$P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \underset{m}{\mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_j), \sigma^2)}$$
$$P(y | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \prod_{j=1}^{m} P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2)$$

- Another motivation: $E[Y(\mathbf{w}, \mathbf{x}_j)] =$

$$y = f + \epsilon$$

$$E(e^{t\epsilon})$$

$$\underset{\epsilon \sim N(0, \sigma^2)}{E[e^{f+\epsilon}]} \qquad \underset{\epsilon \sim N(0, \sigma^2)}{\text{under}}$$

$$Y_j = \omega^T \phi(x_j) + \epsilon$$

$$\underset{N(0,\sigma^2)}{E}\left[ e^{(\omega^T \phi(x_j) + \epsilon)t} \right] \qquad \underset{N(0,\sigma^2)}{E}\left[ e^{\epsilon t} \right]$$

$$= \underset{N(0,\sigma^2)}{E}\left[ e^{\omega^T \phi(x_j) t} \; e^{\epsilon t} \right] = e^{\omega^T \phi(x_j) t} \underset{N(0,\sigma^2)}{E}\left[ e^{\epsilon t} \right]$$

$$\underbrace{\phantom{xxxx}}_{\sigma^2 t^2 / 2}$$

$$\mathcal{N}\left(\omega^T \phi(x_j), \sigma^2\right) \qquad \underset{\text{mapping}}{\text{Given 1-1}} \leftarrow \underbrace{e^{\omega^T \phi(x_j) t + \sigma^2 t^2 / 2}}_{} \; e^{\sigma^2 t^2 / 2}$$

- Linear Model: $Y$ is a linear function of $\phi(\mathbf{x})$, subject to a random noise variable $\varepsilon$ which we believe is 'mostly' around some threshold $\sigma$:

$$Y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$$
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- This allows for the Probabilistic model

$$P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_j), \sigma^2)$$
$$P(y | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \prod_{j=1}^{m} P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2)$$

- Another motivation: $E[Y(\mathbf{w}, \mathbf{x}_j)] = \mathbf{w}^T \phi(\mathbf{x}_j) = \mathbf{w}_0^T + \mathbf{w}_1^T \phi_1(\mathbf{x}_j) + ... + \mathbf{w}_n^T \phi_n(\mathbf{x}_j)$

Sanity check!

$$P(y_j \mid x_j, w) = \mathcal{N}\left(w^T \phi(x_j), \sigma^2\right)$$

Need to estimate "the most representative $w$" for given $D = \{(x_1, y_1), (x_2, y_2) \cdots (x_m, y_m)\}$

most likely

$$L(w \mid D) = P_r\left((x_1, y_1), (x_2, y_2) \cdots (x_m, y_m) \; ; \; w, \sigma^2, \phi(\cdot)\right)$$

- If $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and $y = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$ where $\mathbf{w}, \ \phi(\mathbf{x}) \in \mathbf{R^m}$ then, given dataset $\mathcal{D}$, find the most likely $\hat{\mathbf{w}_{ML}}$

- Recall: $\Pr(y_j | \mathbf{x}_j, \mathbf{w}) = \dfrac{1}{\sqrt{2\pi\sigma^2}} exp\left( \dfrac{(y_j - \mathbf{w}^T\phi(\mathbf{x}_j))^2}{2\sigma^2} \right)$

- From *Probability of data* to *Likelihood of parameters*:

  $\underbrace{\phantom{Likelihood of parameters}}$

  $\boxed{\Pr(\mathcal{D}|\mathbf{w})} = \Pr(\mathbf{y}|\mathbf{x}, \mathbf{w}) =$

  $L(\omega | \mathcal{D})$

  $Pr(y_1 y_2 \cdots y_m | x_1, x_2 \cdots x_m, \phi(\cdot), \omega, \sigma^2) = \prod_i Pr(y_i | x_i)$

  The $(x_1, y_1) \cdots (x_m, y_m)$ collectively influence the fit ($\omega$) **BUT**

  $(x_i, y_i)$   $\bullet(x_j, y_j)$

  Given a $\omega$, prediction $y_j$ for $x_j$ does NOT influence $y_i$ for $x_i$

- If $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and $y = \mathbf{w}^\mathsf{T}\phi(\mathbf{x}) + \epsilon$ where $\mathbf{w}, \ \phi(\mathbf{x}) \in \mathbf{R^m}$ then, given dataset $\mathcal{D}$, find the most likely $\hat{\mathbf{w}_{ML}}$

- Recall: $\Pr(y_j|\mathbf{x}_j, \mathbf{w}) = \dfrac{1}{\sqrt{2\pi\sigma^2}} exp\left(\dfrac{(y_j - \mathbf{w}^\mathsf{T}\phi(\mathbf{x}_j))^2}{2\sigma^2}\right)$

- From *Probability of data* to *Likelihood of parameters*:
  $\Pr(\mathcal{D}|\mathbf{w}) = \Pr(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \displaystyle\prod_{j=1}^{m} \Pr(y_j|\mathbf{x}_j, \mathbf{w}) = \prod_{j=1}^{m} \dfrac{1}{\sqrt{2\pi\sigma^2}} exp\left(\dfrac{(y_j - \mathbf{w}^T\phi(\mathbf{x}_j))^2}{2\sigma^2}\right)$

- Maximum Likelihood Estimate $\hat{\mathbf{w}}_{ML} = \underset{\mathbf{w}}{\mathrm{argmax}} \Pr(\mathcal{D}|\mathbf{w}) = \Pr(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \underline{L(\mathbf{w}|\mathcal{D})}$

# Optimization Trick

- Optimization Trick: Optimal point is invariant under monotonically increasing transformation (such as log )

$$\omega: \quad \text{Tut 1 problem:}$$

$$x^* = \arg\max_x \; \Omega(x) \longrightarrow L(\omega|D)$$

log

Let $\gamma$ be a monontonically increasing fn

Then: $\arg\max_x \; \gamma(\Omega(x)) = x^* \quad (\text{claim})$

# Optimization Trick

- Optimization Trick: Optimal point is invariant under monotonically increasing transformation (such as log )

- $\log L(\mathbf{w}|\mathcal{D}) = LL(\mathbf{w}|\mathcal{D}) = -\frac{m}{2}ln(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{j=1}^{m}(\mathbf{w^T}\phi(\mathbf{x_j}) - \mathbf{y_j})^2$

  For a fixed $\sigma^2$
  $\mathbf{\hat{w}}_{ML} =$

$$\log\left(\prod_{j}^{m}\frac{1}{\sqrt{2\pi\sigma^2}}\,exp\left(-\frac{(w^s\phi(x_j) - y_j)^2}{2\sigma^2}\right)\right) = \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^m exp\left(-\sum_{j}\left[\frac{w^s\phi(x_j) - y_j}{2\sigma^2}\right]^2\right)$$

$$= \frac{-m}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{j=1}(w^s\phi(x_j) - y_j)^2$$

- Optimization Trick: Optimal point is invariant under monotonically increasing transformation (such as log )

- $\log L(\mathbf{w}|\mathcal{D}) = LL(\mathbf{w}|\mathcal{D}) = -\frac{m}{2} ln(\cancel{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{j=1}^{m} (\mathbf{w}^{\mathsf{T}}\phi(\mathbf{x_j}) - \mathbf{y_j})^2$

  *Independent of $\omega$*     *minimize over negated component*

  For a fixed $\sigma^2$
  $\mathbf{w}_{ML}^{\hat{}} = \text{argmax } LL(y_1...y_m|\mathbf{x}_1...\mathbf{x}_m, \mathbf{w}, \sigma^2)$

  $= \underset{\omega}{\text{argmin }} \sum_{j=1}^{m} (\mathbf{w}^T\phi(\mathbf{x}_j) - y_j)^2$

- Note that this is same as the Least square solution!!

  $\hookrightarrow$ with additional power to predict
  $Pr(y_j \mid x_j, \hat{\omega}_{ML}, \sigma^2)$

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori Estimates, Regularization
3. How to minimize the resultant and more complex error functions?
   - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality
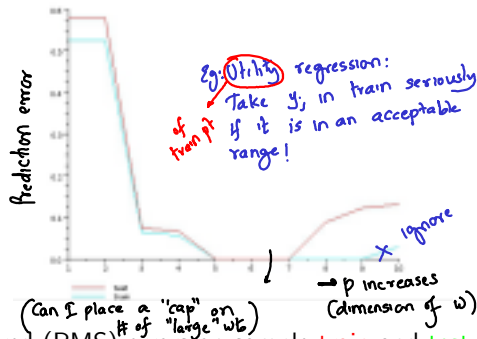
Figure 2: Root Mean Squared (RMS) errors on sample train and test datasets as a function of the degree $t$ of the polynomial being fit

- Too many bends (t=9 onwards) in curve $\equiv$ high values of some $w_i's$. Try plotting values of $w_i$'s using applet at http://mste.illinois.edu/users/exner/java.f/leastsquares/#simulation
- Train and test errors differ significantly

- The Bayesian interpretation of probabilistic estimation is a logical extension that enables reasoning with uncertainty **but in the light of some background belief**
- **Bayesian linear regression**: A Bayesian alternative to **Maximum Likelihood** least squares regression
- Continue with Normally distributed errors
- Model the **w** using a prior distribution and use the posterior over **w** as the result
- Intuitive Prior:

Combining $\epsilon \sim \mathcal{N}(0, \sigma^2)$ with $\omega_i \sim \mathcal{N}(0, \frac{1}{\lambda})$

- The Bayesian interpretation of probabilistic estimation is a logical extension that enables reasoning with uncertainty **but in the light of some background belief**

- **Bayesian linear regression**: A Bayesian alternative to **Maximum Likelihood** least squares regression

- Continue with Normally distributed errors

- Model the **w** using a prior distribution and use the posterior over **w** as the result

- Intuitive Prior: Components of **w** should not become too large!

- Next: Illustration of Bayesian Estimation on a simple Coin-tossing example

Hint from $\epsilon$: $\omega_i \sim \mathcal{N}(0, \frac{1}{\lambda})$

Implicitly putting a cap on [ $\underbrace{\quad}_{\omega}$ ] by restricting 99.7% $\omega_i \in [\pm 3/\sqrt{\lambda}]$

Ideally: $\|\omega\|_0 \leq \theta$

# of non-zero components

Limitation: No probabilistic interpretation

Good news: An efficient algo published in 2015
on minimizing error s.t $\|\omega\|_0 \leq \theta$

[On Iterative Methods for Hard Thresholding,
                    Prateek Jain, MSR I]