

## EECS 3405 Final Review Formulas

Linear Regression: Model	$y = w^T x$
Linear Regression: Loss Function	$E(w) = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$
Linear Regression: Derivative	$\frac{\partial E}{\partial w} = X^T X w - X^T y$
Linear Regression: Gradient for a particular point $(x,y)$	$\begin{aligned} \frac{\partial E}{\partial w} \frac{1}{2} (w^T x - y)^2 &= \frac{1}{2} \cdot 2(w^T x - y) \frac{\partial E}{\partial w} (w^T x - y) \\ &= (w^T x - y)x \\ \nabla E(w) &= (x^T w - y)x \end{aligned}$
Linear Regression: Closed form of derivative by setting to 0	$w^* = (X^T X)^{-1} X^T y$
<b>Stochastic Gradient Descent Method (SGD)</b>	<p>randomly choose <math>\mathbf{x}^{(0)}</math>, and set <math>\eta_0</math>  set <math>n = 0</math></p> <p><b>while</b> not converged <b>do</b></p> <ul style="list-style-type: none"> <li>update: <math>\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \eta_n \nabla f(\mathbf{x}^{(n)})</math></li> <li>adjust: <math>\eta_n \rightarrow \eta_{n+1}</math></li> <li style="text-align: right;"><math>n = n + 1</math></li> </ul> <p><b>end while</b></p>
LDA	<p>Supervised:</p> $\widehat{\mathbf{w}} = \arg \max_w \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$ $S_w^{-1} S_b w = \lambda w$

PCA	Unsupervised: $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (v_i - \hat{v})(v_i - \hat{v})$ $\sigma^2 = w^T S w$
PCA Procedure	<p>1 compute the sample covariance matrix <math>S</math> from training data (<i>measures how features vary</i>)</p> <p>2 calculate the top <math>m</math> eigenvectors of <math>S</math> (<i>Eigenvectors corresponding to largest eigenvalues are the directions where the data has the most variance, Choose the top <math>m</math> of these to reduce the data to <math>m</math> dimensions. Top</i> the largest eigen value)</p> <p>Then, put those selected eigen vectors as a row of a matrix ("A"), size <math>m \times n</math>:</p> <p>3 form <math>A \in \mathbb{R}^{m \times n}</math> with an eigenvector in a row</p> $A = \begin{bmatrix} \quad \hat{w}_1^T \quad \quad \\ \quad \hat{w}_2^T \quad \quad \\ \vdots \\ \quad \hat{w}_m^T \quad \quad \end{bmatrix}_{m \times n}$ <p>4 for any <math>x \in \mathbb{R}^n</math>, map it to <math>y \in \mathbb{R}^m</math> as <math>y = Ax</math>.</p>
Lagrange Theorem	$\nabla f(x, y) = \lambda \nabla g(x, y)$ <p>Set <math>\nabla_x L = 0</math> and <math>\nabla_y L = 0</math></p> <p>Example: <math>f(x, y) = x y</math>, Subject to: <math>g(x, y) = x + y - 10 = 0</math></p> $L(x, y, \lambda) = x y - \lambda(x + y - 10)$ $= x y - \lambda x - \lambda y + \lambda 10$ $\frac{\partial L}{\partial x} = y - \lambda = 0 \quad ; \quad \frac{\partial L}{\partial y} = x - \lambda = 0 \quad ; \quad \frac{\partial L}{\partial \lambda} = x + y - 10 = 0$ <p>Solving this <math>\lambda = x = y</math> and <math>x, y = 5</math> thus: <math>f(5, 5) = 25</math></p>

<b>MCE: Heavyside function which sums the errors</b>	$E_0(w) = \sum_{i=1}^N H(-y_i w^T x_i)$
<b>Sigmoid approximation of MCE</b>	$E_1(w) = \sum_{i=1}^N \sigma(-y_i w^T x_i)$
<b>Gradient of Sigmoid Approximation:</b>	<p>Using derivative of Sigmoid:</p> $\frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$ $\begin{aligned} \frac{\partial}{\partial w} l(-y_i w^T x_i) &= l(-y_i w^T x_i)(1 - l(-y_i w^T x_i)) \cdot \frac{\partial}{\partial w} (-y_i w^T x_i) \\ &= -y_i l(-y_i w^T x_i)(1 - l(-y_i w^T x_i))(x_i) \end{aligned}$ <p>Thus:</p> $\frac{\partial E_1(w)}{\partial w} = - \sum_{i=1}^N y_i l(-y_i w^T x_i)(1 - l(-y_i w^T x_i))(x_i)$
<b>2 ways to represent gradient of MCE:</b>	$\frac{\partial E_1(w)}{\partial w} = - \sum_{i=1}^N y_i l(-y_i w^T x_i)(1 - l(-y_i w^T x_i))(x_i)$ <p>Or, with a simpler gradient:</p> $\frac{\partial E_2(w)}{\partial w} = \sum_{i=1}^N y_i l(-y_i w'^T x_i)(1 - l(y_i w'^T x_i))(x_i)$
<b>The average MCE loss function, show also the vectorized form</b>	<p>Both are average loss (times 1/n) over all samples:</p> $\frac{1}{n} \sum_{i=1}^N l(-y_i w^T x_i)$ $\frac{1}{n} \sigma(-y \odot (Xw))$

<b>Logistic Regression</b>	$w_{LR} = \arg \max_w L(w) = \arg \max_w \prod_i^N \sigma(y_i w^T x_i)$ Where $\sigma = 1$
<b>Logistic Regression in natural logarithmic form and how it is derived:</b>	$\ln(L(w)) = \ln\left(\prod_i^N \sigma(y_i w^T x_i)\right)$ Log property: $\ln(ab)=\ln a + \ln b$ $\ln(a \cdot b \cdot c \dots) = \ln a + \ln b + \ln c + \dots$ $\ln(L(w)) = \sum_{i=1}^N \ln(\sigma(y_i w^T x_i))$
<b>Derivative of Logistic Regression in natural log form:</b>	Taking the gradient wrt w: $\frac{\partial \ln(L(w))}{\partial w} = \sum_{i=1}^N \frac{\partial}{\partial w} \ln(\sigma(y_i w^T x_i))$ $= \sum_{i=1}^N \frac{1}{\sigma(y_i w^T x_i)} \cdot \frac{\partial}{\partial w} \sigma(y_i w^T x_i)$ Derivative of sigmoid function (plus chain rule): $= \sum_{i=1}^N \frac{1}{\sigma(y_i w^T x_i)} \cdot \sigma(y_i w^T x_i) (1 - \sigma(y_i w^T x_i)) (y_i x_i)$ $= \sum_{i=1}^N \frac{1}{1} \cdot (1 - \sigma(y_i w^T x_i)) (y_i x_i)$ $= \sum_{i=1}^N (1 - \sigma(y_i w^T x_i)) (y_i x_i)$

	<p>Note: <math>y_i</math> is a scalar, thus acts like a coefficient:</p> $\frac{\partial \ln(L(w))}{\partial w} = \sum_{i=1}^N y_i(1 - \sigma(y_i w^T x_i))x_i$
<b>How to represent the gradient for a particular mini batch of Logistic Regression (also with element-wise multiplication)</b>	$\begin{aligned}\frac{\partial Q(w; B)}{\partial B} &= \frac{1}{B} \sum_{i=1}^B y_i(1 - \sigma(y_i w^T x_i))x_i \\ &= \frac{1}{B} X^\top [y \odot l(y \odot (Xw)) - y]\end{aligned}$
<b>Logistic Regression Loss Function</b>	$-\ln(\sigma(y \odot (Xw))) \quad (\in \mathbb{R}^N)$
<i>SVM0 Formulation:</i>	$\frac{\max_{\varepsilon, w, b} \gamma}{\ w\ } \geq \gamma, \forall i \in \{1, 2, \dots, N\}$
<b>SVM1: Primal Problem</b>	Primal Problem: $\min_{w,b} \frac{1}{2} w^T w \quad \text{subject to} \quad y_i(w^T x_i + b) \geq 1, \quad \forall i \in \{1, 2, \dots, N\}$
<b>How to represent SVM Primal Problem as a Lagrangian</b>	$\mathcal{L}(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i (y_i(w^T x_i + b) - 1)$ $\mathcal{L}(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i y_i w^T x_i + \alpha_i y_i b - \alpha_i$

	$\mathcal{L}(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i y_i w^T x_i - b \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i$
<b>SVM2: Dual Problem</b>	$\max_{\alpha_i \geq 0} \left[ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \right]$
<b>Quadratic Programming Formulation</b>	$\max_{\alpha} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T Q \alpha$ <p><b>subject to:</b></p> $y^T \alpha = 0$ $\alpha \geq 0$
<b>Define each matrix for Quadratic Programming:</b>	$\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix}_{N \times 1} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1} \quad \mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{N \times 1}$ $\mathbf{Q} = \begin{bmatrix} Q_{ij} \end{bmatrix}_{N \times N} = \begin{bmatrix} \mathbf{y} \mathbf{y}^T \end{bmatrix}_{N \times N} \odot \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \cdots & \mathbf{x}_1^T \mathbf{x}_N \\ \vdots & \mathbf{x}_i^T \mathbf{x}_j & \vdots \\ \mathbf{x}_N^T \mathbf{x}_1 & \cdots & \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix}_{N \times N}$ <p>where <math>\odot</math> denotes element-wise multiplication, and:</p> $Q_{ij} = (y_i y_j) \odot (x_i^T x_j)$

<b>Soft SVM Problem</b>	$\min_{w,b,\xi_i} \left( \frac{1}{2} \right) w^T w + C \sum_{i=1}^N \xi_i$ <p style="text-align: center;"><b>subject to:</b></p> $y_i(w^T x_i + b) \geq 1 - \xi_i, \forall i \in \{1, 2, \dots, N\}$ $\xi_i \geq 0, \forall i \in \{1, 2, \dots, N\}$
<b>Soft SVM: Dual Problem + How is it different?</b>	$\max_{\alpha} 1^T \alpha - \frac{1}{2} \alpha^T Q \alpha$ <p style="text-align: center;"><b>subject to:</b></p> $y^T \alpha = 0$ $0 \leq \alpha \leq C$ <p>the <math>\alpha</math>'s must be greater than 0, but less than C. That C is the hyperparameter for the maximum allowance of error.</p>
<b>L<sub>2</sub> norm</b>	<b>L<sub>2</sub> norm:</b> $\ w\ _2 = \sqrt{( w_1 ^2 + \dots +  w_n ^2)}$
<b>Ridge Regression Objective Function (also in elementwise multiplication)</b>	$w_{ridge}^* = \operatorname{argmin}_w \left[ \sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \cdot \ w\ _2^2 \right]$ $w^* = \operatorname{argmin}_w \left[ \sum_{i=1}^N \left( \sum_{j=1}^d w_j x_{ij} - y_i \right)^2 + \lambda \sum_{j=1}^d (w_j \odot w_j) \right]$
<b>How to derive the Ridge Regression Closed form Solution</b>	$\sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \cdot \ w\ _2^2$ $L_{ridge}(w) = \ Xw - y\ ^2 + \lambda w^T w$ $= (Xw - y)^T (Xw - y) + \lambda w^T w$ $= (Xw - y)^T (Xw - y) + \lambda w^T w$ $= w^T X^T X w - 2Xw y^T + y^T y + \lambda w^T w$ $\nabla L_{ridge}(w) = 2X^T X w - 2X^T y + 2\lambda w$

	<p>Set gradient to 0:</p> $2X^T X w - 2X^T y + 2\lambda w = 0$ $X^T X w - X^T y + \lambda w = 0$ $(X^T X + \lambda I)w = X^T y$ $w_{ridge}^* = (X^T X + \lambda I)^{-1} X^T y$
<b>LASSO Objective Function</b>	$w_{lasso}^* = \underset{w}{\operatorname{argmin}} \left[ \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \cdot \ w\ _1 \right]$
<b>How to derive LASSO Gradient (component wise)</b>	<p><math>\text{sign}(w_j) = \begin{cases} +1 &amp; \text{if } w_j &gt; 0 \\ -1 &amp; \text{if } w_j &lt; 0 \\ 0 &amp; \text{if } w_j = 0 \text{ (convention for gradient descent)} \end{cases}</math></p> $\begin{aligned} Q_{lasso}(w) &= \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \cdot \ w\ _1 \\ &= \sum_{i=1}^N (w^T x_i - y_i) \cdot x_i + \lambda \cdot \text{sign}(w) \\ &= \sum_{i=1}^N x_i (w^T x_i) - \sum_{i=1}^N y_i x_i + \lambda \cdot \text{sign}(w) \\ &= \sum_{i=1}^N x_i (x_i^T w) - \sum_{i=1}^N y_i x_i + \lambda \cdot \text{sign}(w) \\ &= \left( \sum_{i=1}^N x_i x_i^T \right) w - \sum_{i=1}^N y_i x_i + \lambda \cdot \text{sign}(w) \end{aligned}$

<b>Loss function for matrix factorization</b>	$Q(U, V) = \sum_{(i,j) \in \Omega} (x_{ij} - u_i^T v_j)^2 + \lambda_1 \sum_{i=1}^n \ u_i\ _2^2 + \lambda_2 \sum_{j=1}^m \ v_j\ _2^2$
<b>Neural Network Softmax function</b>	$y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$
<i>Forward Pass of a Fully-Connected DNN</i>	<p>1. For the input layer: <math>z_0 = x</math></p> <p>2. For each hidden layer <math>l = 1, 2, \dots, L-1</math>:</p> <p><i>The output from previous layer (<math>a_l</math>) equals the non-linear activation plus bias:</i></p> $a_l = W^{(l)} z_{l-1} + b^l$ $z_l = \text{ReLU}(a_l)$ <p>3. For the output layer:</p> $a_L = W^{(L)} z_{L-1} + b^{(L)}$ $y = z_L = \text{softmax}(a_L)$
<i>Derivatives for different automatic differentiation</i>  (full connection, nonlinear activation)	
Why is $x^T Q x$ considered quadratic? (where $x$ is a column vector and $Q$ is a matrix)	$x^T Q x = \sum_{i=1}^n \sum_{j=1}^n Q_{ij} x_i x_j$ <p><b>Every term is of degree 2</b></p>

*How to solve a max / min problem based on a matrix*

Q:

- Solve  $\nabla f = 0 \rightarrow$  this gives critical point
- Check Hessian =  $2Q$
- If  $Q > 0 \rightarrow$  it's a **minimum**
- If  $Q < 0 \rightarrow$  it's a **maximum**

If Q indefinite  $\rightarrow$  **saddle**