

Properties of trace:

① We will first show the cyclic property of Trace.

$$\text{Tr}(AB) = \text{Tr}(BA)$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times n}$.

Let $A = \begin{bmatrix} -a_1^T & - \\ -a_2^T & - \\ \vdots & \\ -a_n^T & \end{bmatrix}$, $B = \begin{bmatrix} 1 & 1 & \dots & 1 \\ b_1 & b_2 & \dots & b_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$

Then,
 $AB = \begin{bmatrix} a_1^T b_1 & a_1^T b_2 & \dots & a_1^T b_n \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ a_n^T b_1 & a_n^T b_2 & \dots & a_n^T b_n \end{bmatrix}$

so, $\text{tr}(AB) = \sum_{i=1}^n a_i^T b_i$

Now,

$$BA = b_1 a_1^T + b_2 a_2^T + \dots + b_n a_n^T$$

Since trace is a linear operator, so

$$\text{tr}(BA) = \sum_{i=1}^n \text{tr}(b_i a_i^T)$$

Now,

$$b_i a_i^T = \begin{bmatrix} b_{1i} \\ b_{2i} \\ \vdots \\ b_{ni} \end{bmatrix} [a_{i1} \ a_{i2} \ \dots \ a_{in}]$$

So,

$$\text{tr}(b_i a_i^T) = \sum_{m=1}^n a_{im} b_{mi}$$

$$\Rightarrow \text{tr}(b_i a_i^T) = a_i^T b_i$$

$$\therefore \text{tr}(BA) = \sum_{i=1}^n a_i^T b_i = \text{tr}(AB)$$

② By definition of trace, it's trivial that

$$\text{tr}(A) = \text{tr}(A^T)$$

③ Let's show,

$$\text{tr}(A^T B) = \sum_i \sum_j a_{ij} b_{ij}$$

Let $A = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix}$, $B = \begin{bmatrix} b_1 & b_2 & \dots & b_n \end{bmatrix}$

So, $A^T B = \begin{bmatrix} -a_1^T - \\ -a_2^T - \\ \vdots \\ -a_n^T - \end{bmatrix} \begin{bmatrix} b_1 & b_2 & \dots & b_n \end{bmatrix}$

Now, $\text{tr}(A^T B) = \sum_{i=1}^n a_i^T b_i$

where $a_i = \begin{bmatrix} a_{1i} \\ a_{2i} \\ \vdots \\ a_{ni} \end{bmatrix}$, $b_i = \begin{bmatrix} b_{1i} \\ \vdots \\ b_{ni} \end{bmatrix}$

from above we can see that we are multiplying A and B elementwise and then summing up to get the trace of $A^T B$.

$$\text{So, } \text{tr}(A^T B) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ij}$$

chain rule of derivatives:

Let $y = f(x)$ for $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$
 and $z = g(y)$ for $g: \mathbb{R}^n \rightarrow \mathbb{R}^p$. Then

we have,

$$\nabla_x z = \nabla_x y \nabla_y z$$

Let's first do a dimensionality check on the chain rule. We know

that using denominator layout

$$\nabla_x z \in \mathbb{R}^{m \times p}, \nabla_{xy} \in \mathbb{R}^{m \times n}, \nabla_y z \in \mathbb{R}^{n \times p}$$

Hence,

$$\nabla_x y \nabla_y z \in \mathbb{R}^{m \times p}$$

Therefore the dimensions in the chain

rule is consistent. Now let's do an

intuition check. From the derivative

of a vector with respect to a vector

we know,

$$\Delta z \approx (\nabla_x z)^T \Delta x$$

The chain rule is intuitive since

$$\Delta y \approx (\nabla_x y)^T \Delta x$$

$$\Delta z \approx (\nabla_y z)^T \Delta y$$

composing we get,

$$\Delta z \approx (\nabla_y z)^T (\nabla_x y)^T \Delta x$$

Hence,

$$(\nabla_x z)^T = (\nabla_y z)^T (\nabla_x y)^T$$

Practice problem on derivatives:

① Let

$$A = \begin{bmatrix} -a_1^T - \\ -a_2^T - \\ \vdots \\ -a_n^T - \end{bmatrix}, \quad B = \begin{bmatrix} b_1 & b_2 & \dots & b_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

Then,

$$\text{Tr}(AB) = a_1^T b_1 + a_2^T b_2 + \dots + a_n^T b_n$$

Now, using the definition

$$\begin{aligned} \nabla_A \text{Tr}(AB) &= \frac{\partial \text{Tr}(AB)}{\partial a_{ij}} \\ &= \frac{\partial}{\partial a_{ij}} \left[\sum_{m=1}^n a_{im} b_{mj} + \dots + \sum_{m=1}^n a_{nm} b_{mn} \right] \\ &= b_{ji} \\ \therefore \nabla_A \text{Tr}(AB) &= B^T \end{aligned}$$

$$\textcircled{2} \quad \nabla_A (x^T A x)$$

since $\text{Tr}(X^T A X) = X^T A X$, so

$$= \nabla_A [\text{Tr}(X^T A X)]$$

Using the cyclic property, we have

$$= \nabla_A [\text{Tr}(A X X^T)]$$

Using the result from \textcircled{1},

$$= X X^T$$

$$\therefore \nabla_A X^T A X = X X^T$$

③ We want to compute

$$\nabla_z (x-z)^T \Sigma^{-1} (x-z)$$

Let $y = f(z) = x-z$
 $r = g(y) = y^T \Sigma^{-1} y$

Then using chain rule of derivatives,

$$\nabla_z r = \nabla_z y \nabla_y r$$

Now,
 $\nabla_z y = -I \in \mathbb{R}^{n \times n}$

From lecture we know,

$$\nabla_y (y^T \Sigma^{-1} y) = [(\Sigma^{-1} + [\Sigma^{-1}]^T) y]$$

since Σ^{-1} is symmetric, so

$$\nabla_y (y^T \Sigma^{-1} y) = 2 \Sigma^{-1} y$$

$$\therefore \nabla_{\theta} (\theta^T (\theta - \bar{\theta})) = -2 \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)T}$$

Regularized least squares:

In RLS, we want to solve the following optimization problem

$$\arg \min_{\theta} \frac{1}{2} \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

From lecture we know that

$$\frac{1}{2} \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)})^2 = \frac{1}{2} (\gamma - X\theta)^T (\gamma - X\theta)$$

where, $\gamma = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times 1}$

$$X = \begin{bmatrix} \hat{x}^{(1)T} \\ \vdots \\ \hat{x}^{(N)T} \end{bmatrix} \in \mathbb{R}^{N \times 2}$$

Then the optimization problem becomes

$$\arg \min_{\theta} \frac{1}{2} (Y - X\theta)^T (Y - X\theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

Since the cost function is convex, so we can solve for θ^* by setting the derivative w.r.t θ to zero.

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{2} [Y^T Y - Y^T X\theta - \theta^T X^T Y + \theta^T X^T X\theta] \\ &\quad + \frac{\lambda}{2} \theta^T \theta \end{aligned}$$

$$= \frac{1}{2} [Y^T Y - 2Y^T X\theta + \theta^T X^T X\theta] + \frac{\lambda}{2} \theta^T \theta$$

Now,

$$\nabla_{\theta} \mathcal{L} = -\nabla_{\theta} [Y^T X\theta] + \frac{1}{2} \nabla_{\theta} [\theta^T X^T X\theta] + \frac{\lambda}{2} \nabla_{\theta} [\theta^T \theta]$$

$$\begin{aligned} \nabla_{\theta} \mathcal{L} &= -\nabla_{\theta} [Y^T X\theta] + \frac{1}{2} \nabla_{\theta} [\theta^T (X^T X + \lambda I)\theta] \\ &= -\nabla_{\theta} [Y^T X\theta] + \frac{1}{2} \nabla_{\theta} [\theta^T (X^T X + \lambda I)\theta] \end{aligned}$$

$$= -X^T Y + (X^T X + \lambda I)\theta$$

Setting $\nabla_{\theta} \mathcal{L} = 0$ and solving for θ we get

$$\theta^* = (X^T X + \lambda I)^{-1} X^T Y$$

1.5.4

$$f = \|A + \lambda B\|_F^2$$

Recall, for a matrix X the frobenius norm is defined as follows

$$\|X\|_F^2 = \text{tr}(X^T X)$$

using the above definition,

$$\begin{aligned} f &= \text{tr}[(A + \lambda B)^T (A + \lambda B)] \\ &= \text{tr}[A^T A] + \lambda \text{tr}[A^T B] + \lambda \text{tr}[B^T A] \\ &\quad + \lambda^2 \text{tr}[B^T B] \end{aligned}$$

since B is an orthogonal matrix, so

$$B^T B = I$$

Hence,

$$\begin{aligned} f &= \text{tr}[A^T A] + \lambda \text{tr}[A^T B] + \lambda \text{tr}[B^T A] \\ &\quad + \lambda^2 n \end{aligned}$$

Dropping terms that has no B dependence, so

$$f = \lambda \text{Tr}[A^T B] + \lambda \text{Tr}[B^T A]$$

Since, $\text{Tr}[A^T B] = \text{Tr}[B^T A]$

so,
 $f = 2\lambda \text{Tr}[A^T B]$

Hence,

$$\nabla_B f = 2\lambda \nabla_B \text{Tr}[A^T B]$$
$$= 2\lambda A \quad (\text{Using (103)} \\ \text{from matrix cookbook})$$

Word translation with linear models:

There are several ways to convert words into their vector representations. Some of the popular word embedding models are:

(i) Word 2 Vec

(ii) Glove

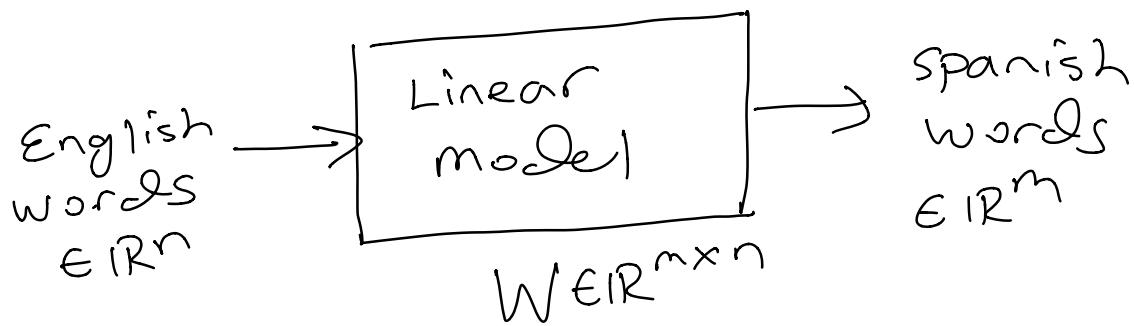
(iii) BERT

:

:

We won't get into the details of word embedding models in this discussion, but by the end of this quarter you will have all the tools necessary to

Understand them.



a)
$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^K \|b^{(i)} - Wa^{(i)}\|_2^2$$

The loss function is measuring the L_2 distance between the Spanish words in \mathbb{R}^m and the linearly mapped English words in \mathbb{R}^m .

The goal of the optimization is to learn a linear mapping W such that the mapped English words are as close as possible to their Spanish counterparts in the Euclidean norm (L_2) sense.

If there exists a linear mapping between English and Spanish words, then the linear transformation W , will map the English word 'dog' to the Spanish word 'perro'.
to the Spanish word 'perro'.

$$b) L(w) = \frac{1}{2} \sum_{i=1}^k \|b^{(i)} - wa^{(i)}\|_2^2$$

Recall that for a matrix $X \in \mathbb{R}^{m \times n}$

$$\|X\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n (x_{ij})^2$$

$$= \sum_{j=1}^n \|X(:, j)\|_2^2$$

So if we define the following
matrices

$$B = \begin{bmatrix} b^{(1)} & b^{(2)} & \dots & b^{(K)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{m \times K}$$

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(K)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times K}$$

Then,
 $L(w) = \frac{1}{2} \|B - wA\|_F^2$
 using the trace definition of
 frobenius norm,

$$L(w) = \frac{1}{2} \text{Tr} \left[(B - wA)^T (B - wA) \right]$$

$$= \frac{1}{2} \text{Tr} \left[(B^T - A^T w^T) (B - wA) \right]$$

$$= \frac{1}{2} \text{Tr} \left[B^T B - B^T w A - A^T w^T B + A^T w^T w A \right]$$

Dropping the terms that has
 no w dependence and observing

$$\begin{aligned} \text{Tr}[B^T w A] &= \text{Tr}[(B^T w A)^T] \\ &= \text{Tr}[A^T w^T B] \end{aligned}$$

we have

$$L(w) = -\text{Tr}[A^T w^T B] + \frac{1}{2} \text{Tr}[A^T w^T w A]$$

Now by cyclic property of trace

$$\begin{aligned} & \text{Tr}[A^T w^T w A] \\ &= \text{Tr}[w^T w A A^T] \end{aligned}$$

Hence,

$$L(w) = -\text{Tr}[A^T w^T B] + \frac{1}{2} \text{Tr}[w^T w A A^T]$$

Now,

$$\nabla_w \text{Tr}[A^T w^T B] = B A^T \quad (\text{(102) in matrix cookbook})$$

$$\nabla_w \text{Tr}[w^T w A A^T] = 2 w A A^T \quad (\text{(113) in matrix cookbook})$$

So,

$$\nabla_w L(w) = -BA^T + WAA^T$$