# BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda \|w\|^2$$

$\hookrightarrow$ 1. expand

2. take derivative
   w.r.t to w

3. set it to 0

$$\|$$
$$\vee$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

## BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda\|w\|^2$$

$\llcorner$ 1. expand

2. take derivative w.r.t to w

3. set it to 0

$$\Updownarrow$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

## DUAL RIDGE

- approach from mathematical programming (field of OR/LP)

## BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda \|w\|^2$$

$\llcorner$ 1. expand

   2. take derivative
      w.r.t to w

   3. set it to 0

$$\|$$
$$\downarrow$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

## DUAL RIDGE

- approach from mathe-
matical programming
(field of OR/LP)

objective: $\arg\min_{w, e} \|e\|^2 + \lambda \|w\|^2$

s.t. $e = X^T w - y$

$\llcorner$ variables $e, w$

  - coefficients $X, y$

# BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda \|w\|^2$$

$\mathrel{L_{>}}$ 1. expand

2. take derivative w.r.t to w

3. set it to 0

$$\|$$
$$\Downarrow$$

$$w = (X^TX - \lambda I)^{-1} X^T y$$

# DUAL RIDGE

- approach from mathematical programming (field of OR/LP)

objective: $\arg\min_{w, e} \|e\|^2 + \lambda \|w\|^2$

s.t. $e = X^T w - y$ $\longleftarrow$

$\mathrel{L_{>}}$ variables $e, w$
- coefficients $X, y$

$\Rightarrow$ reformulate into Lagrangian

key assumption: $X$ is centred at mean $0$, thus **NO** intercept

# BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda \|w\|^2$$

$\downarrow$ 1. expand

2. take derivative w.r.t to $w$

3. set it to 0

$$\|$$
$$\downarrow$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

# DUAL RIDGE

- approach from mathematical programming (field of OR/LP)

objective: $\arg\min_{w, e} \|e\|^2 + \lambda \|w\|^2$

s.t. $e = X^T w - y$ $\leftarrow$ 

key assumption: $X$ is centred at mean 0, thus **NO** intercept

$\downarrow$ variables $e, w$
- coefficients $X, y$

$\Rightarrow$ reformulate into Lagrangian

lagrange multipliers
$\downarrow$

$$L(e, w, \alpha) = \|e\|^2 + \lambda \|w\|^2 + \alpha^T (e - X^T w + y)$$

# BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda \|w\|^2$$

↳ 1. expand

2. take derivative w.r.t to $w$

3. set it to 0

$$\|$$
$$\downarrow$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

# DUAL RIDGE

- approach from mathematical programming (field of OR/LP)

objective: $\arg\min_{w, e} \|e\|^2 + \lambda \|w\|^2$

s.t. $e = X^T w - y$ ← key assumption: $X$ is centred at mean 0, thus **NO** intercept

↳ variables $e, w$
  - coefficients $X, y$

⟹ reformulate into Lagrangian

$$L(e, w, \alpha) = \|e\|^2 + \lambda \|w\|^2 + \alpha^T (e - X^T w + y)$$

→ derivate!

$$\frac{\partial L}{\partial e} = 2e + \alpha \qquad \left.\right\} \text{set to } \emptyset$$

$$\frac{\partial L}{\partial w} = 2\lambda w - X\alpha$$

lagrange multipliers
↓

## BASIC RIDGE

$$\min_{w} \|Xw - y\|^2 + \lambda\|w\|^2$$

$\llcorner$ 1. expand

2. take derivative w.r.t to $w$

3. set it to $0$

$$\|$$

$$w = (X^T X - \lambda I)^{-1} X^T y$$

## DUAL RIDGE

- approach from mathematical programming (field of OR/LP)

objective: $\arg\min_{w, e} \|e\|^2 + \lambda\|w\|^2$

s.t. $e = X^T w - y$

$\llcorner$ • variables $e, w$
- coefficients $X, y$

$\Rightarrow$ reformulate into Lagrangian

$$L(e, w, \alpha) = \|e\|^2 + \lambda\|w\|^2 + \alpha^T(e - X^T w + y)$$

key assumption: $X$ is centred at mean $0$, thus **NO** intercept

lagrange multipliers $\downarrow$

$\rightarrow$ derivate!

$$\left.\begin{array}{l} \dfrac{\partial L}{\partial e} = 2e + \alpha \\[2mm] \dfrac{\partial L}{\partial w} = 2\lambda w - X\alpha \end{array}\right\} \begin{array}{l} \text{set} \\ \text{to } \emptyset \end{array}$$

$$\boxed{\begin{array}{c} e = -\dfrac{1}{2}\alpha^T \\[2mm] w = \dfrac{1}{2\lambda} X\alpha \end{array}}$$

optimal values of $w, e$

$\llcorner$ let's substitute them in!

# Dual Ridge (continued)

$$L(e(\alpha), w(\alpha), \alpha) =$$

$$= \frac{1}{4}\|\alpha\|^2 + \frac{\lambda}{4\lambda^2}\|X\alpha\|^2 + \alpha^T\left(-\frac{1}{2}\alpha - \frac{1}{2\lambda}X^TX\alpha + y\right)$$

$$= -\frac{1}{4}\|\alpha\|^2 - \frac{1}{4\lambda}\|X\alpha\|^2 + \alpha^T y$$

# Dual Ridge (continued)

$$L(e(\alpha), w(\alpha), \alpha) =$$

$$= \frac{1}{4}\|\alpha\|^2 + \frac{\lambda}{4\lambda^2}\|X\alpha\|^2 + \alpha^T\left(-\frac{1}{2}\alpha - \frac{1}{2\lambda}X^TX\alpha + y\right)$$

$$= -\frac{1}{4}\|\alpha\|^2 - \frac{1}{4\lambda}\|X\alpha\|^2 + \alpha^T y$$

$\quad\quad\hookrightarrow X$: known

$\quad\quad\quad\quad y$: known

$\quad\quad\quad\quad \lambda$: we choose

$\quad\quad\quad\quad \alpha$: <u>unknown!</u>

$\quad\quad\quad\quad\quad\quad \underset{\Vdash}{} \underline{\text{DUAL}}: \quad \arg\min_{\alpha}\left[-\frac{1}{4}\alpha^T\left(I + \frac{1}{\lambda}X^TX\right)\alpha + \alpha^T y\right] \quad /\cdot(-\lambda)$

$$g(\alpha) = \arg\min_{\alpha}\left[\frac{1}{4}\lambda\,\alpha^T(X^TX - \lambda I)\alpha - \lambda\alpha^T y\right]$$

# Dual Ridge (continued)

$$L(e(\lambda), w(\lambda), \lambda) =$$

$$= \frac{1}{4}\|\lambda\|^2 + \frac{\lambda}{4\lambda^2}\|X\lambda\|^2 + \lambda^T\left(-\frac{1}{2}\lambda - \frac{1}{2\lambda}X^TX\lambda + y\right)$$

$$= -\frac{1}{4}\|\lambda\|^2 - \frac{1}{4\lambda}\|X\lambda\|^2 + \lambda^T y$$

$\llcorner$ $X$ : known

  $y$ : known

  $\lambda$ : we choose

  $\lambda$ : __unknown!__

$\lfloor\!\!\!\!L$> __DUAL__ : $\arg\min\limits_{\lambda}\left[-\frac{1}{4}\lambda^T\left(I + \frac{1}{\lambda}X^TX\right)\lambda + \lambda^T y\right]$

$$g(\lambda) = \arg\min\limits_{\lambda}\left[\frac{1}{4}\lambda\,\lambda^T(X^TX - \lambda I)\lambda - \lambda\lambda^T y\right]$$

$\frac{\partial g}{\partial \lambda} = \frac{1}{2}\lambda^T(X^TX - I\lambda) - \lambda y \qquad = 0$

$$\lambda = 2\lambda(X^TX - I\lambda)^{-1}y$$

$\lfloor\!\!\!\!L$> $w = X(X^TX - I\lambda)^{-1}y$

!EQUIVALENT TO PRIMAL!

$/ \cdot (-\lambda)$

# Dual Ridge (continued)

$$L(e(\lambda), w(\lambda), \lambda) =$$

$$= \frac{1}{4} \|\lambda\|^2 + \frac{\lambda}{4\lambda^2} \|X\lambda\|^2 + \lambda^T \left(-\frac{1}{2}\lambda - \frac{1}{2\lambda}X^T X \lambda + y\right)$$

$$= -\frac{1}{4}\|\lambda\|^2 - \frac{1}{4\lambda}\|X\lambda\|^2 + \lambda^T y$$

$$\frac{\partial g}{\partial \lambda} = \frac{1}{2}\lambda^T(X^T X - I\lambda) - \lambda y \bigg/ = 0$$

$$\lambda = 2\lambda(X^T X - I\lambda)^{-1} y$$

$$\Lsh w = X(X^T X - I\lambda)^{-1} y$$

!EQUIVALENT TO PRIMAL!

$\Lsh$ $X$: known

$\quad y$: known

$\quad \lambda$: we choose

$\quad \lambda$: underline{unknown!}

$\Lsh$ DUAL : $\arg\min\limits_{\lambda}\left[-\frac{1}{4}\lambda^T\left(I + \frac{1}{\lambda}X^T X\right)\lambda + \lambda^T y\right]$ $\bigg/ \times(-\lambda)$

$$g(\lambda) = \arg\min\limits_{\lambda}\left[\frac{1}{4\lambda}\lambda^T(X^T X - \lambda I)\lambda - \lambda \lambda^T y\right]$$

# Kernel Regression

$\hookrightarrow$ what is a kernel?

# Kernel Regression

↳ what is a kernel?

  ↳ implicit mapping into a higher dimensional space that gives us pairwise distances of the mapped data points but **NOT** their actual coordinates.

# Kernel Regression

↳ what is a kernel?

↳ implicit mapping into a higher dimensional space that gives us **pairwise distances** of the mapped data points, but <u>NOT</u> their actual coordinates.

$X^T X$ → <u>cov. matrix</u>
aka
<u>dot product matrix</u>
aka
<u>matrix of pair-wise distances</u>

# Kernel Regression

$\rightarrow$ what is a kernel?

$\rightarrow$ implicit mapping into a higher dimensional space that gives us **pairwise distances** of the mapped data points, but __NOT__ their actual coordinates.

$$X^T X \rightarrow \underline{cov. \ matrix}$$

aka
dot product matrix
aka
matrix of pair-wise distances

$\rightarrow$ since dual ridge only needs $X^T X$ it provides a perfect opportunity to replace $X^T X = K = \Phi(x) \cdot \Phi(x)$

# Understanding Why We need $\bar{X} = 0$

$\downarrow$ usual ridge: $f(w_0, w) = \sum_{i=1}^{n} (x_i^T w + w_0 - y_i)^2 + \lambda \| w \|^2$

$\qquad \downarrow \dfrac{\partial f}{\partial w_0} = 2 \sum_{i=1}^{n} x_i^T w + w_0 - y_i = 0$

# Understanding Why We need $\bar{X} = 0$

$\rightarrow$ usual ridge: $\quad f(w_0, w) = \sum_{i=1}^{n} (x_i^T w + w_0 - y_i)^2 + \lambda \|w\|^2$

$\qquad\qquad\qquad \rightarrow \dfrac{\partial f}{\partial w_0} = 2 \sum_{i=1}^{n} x_i^T w + w_0 - y_i \quad = \quad 0$

$\qquad\qquad\qquad\qquad\qquad n w_0 = \sum_i y_i - \sum_i x_i^T w$

$\qquad\qquad\qquad\qquad\qquad w_0 = \bar{y} - \bar{X}^T w$

# Understanding why we need $\bar{X} = 0$

$\hookrightarrow$ usual ridge: $\quad f(w_0, w) = \sum\limits_{i=1}^{n} (x_i^T w + w_0 - y_i)^2 + \lambda \|w\|^2$

$\qquad\qquad\qquad\qquad \hookrightarrow \dfrac{\partial f}{\partial w_0} = 2 \sum\limits_{i=1}^{n} x_i^T w + w_0 - y_i = 0$

$\qquad\qquad\qquad\qquad\qquad\qquad n w_0 = \sum\limits_i y_i - \sum\limits_i x_i^T w$

$f(w, w_0(w)) = \sum\limits_{i=1}^{n} \left[ (x_i - \bar{x})^T w - (y_i - \bar{y}) \right]^2 \Longleftarrow \qquad w_0 = \bar{y} - \bar{x}^T w$
$\qquad\qquad\qquad\quad + \lambda \|w\|^2$

# Understanding Why We need $\bar{X} = 0$

$\hookrightarrow$ usual ridge: $f(w_0, w) = \sum_{i=1}^{n} (x_i^T w + w_0 - y_i)^2 + \lambda \|w\|^2$

$\qquad \hookrightarrow \dfrac{\partial f}{\partial w_0} = 2 \sum_{i=1}^{n} x_i^T w + w_0 - y_i = 0$

$\qquad\qquad\qquad n w_0 = \sum_i y_i - \sum_i x_i^T w$

$f(w, w_0(w)) = \sum_{i=1}^{n} \left[ (x_i - \bar{x})^T w - (y_i - \bar{y}) \right]^2 \longleftarrow \qquad w_0 = \bar{y} - \bar{x}^T w$

$\qquad\qquad\qquad + \lambda \|w\|^2$

$\underbrace{\qquad\qquad\qquad\qquad}$

if $\bar{x}$ & $\bar{y}$ had 0
mean life would
be simpler

# Understanding Why We need $\bar{X} = 0$

$\rightarrow$ usual ridge: $f(w_0, w) = \sum_{i=1}^{n} (x_i^T w + w_0 - y_i)^2 + \lambda \|w\|^2$

$\qquad \rightarrow \frac{\partial f}{\partial w_0} = 2 \sum_{i=1}^{n} x_i^T w + w_0 - y_i = 0$

$\qquad\qquad\qquad n w_0 = \sum_i y_i - \sum_i x_i^T w$

$f(w, w_0(w)) = \sum_{i=1}^{n} [(x_i - \bar{x})^T w - (y_i - \bar{y})]^2 \longleftarrow \qquad w_0 = \bar{y} - \bar{x}^T w$
$\qquad\qquad + \lambda \|w\|^2$

$\underbrace{\qquad\qquad\qquad\qquad}$

if $\bar{x}$ & $\bar{y}$ had 0
mean life would
be simpler $\implies$

thus we center using
$JX = (I - \frac{1}{n} \mathbb{1}\mathbb{1}^T) X$

$\qquad \rightarrow x_i - \frac{1}{n}\bar{x} \qquad \forall i$