

MODEL FITTING II: RIDGE & TIKHONOV

Prof. Alexander Huth
2.20.2020

RECAP

- * Linear regression!
- * Ordinary least squares (OLS)
- * Regularized regression
 - * Priors on weights
 - * Penalties on weights
 - * Ad hoc metrics (early stopping)

RECAP

- * L2 penalty = Gaussian prior = ridge
- * L1 penalty = Laplacian prior = LASSO

TODAY

- * Analytic solutions to L2-regularized regression problems:
 - * Ridge regression
 - * Tikhonov regression

RIDGE REGRESSION

- * Multivariate normal (MVN) prior on β
- * L2 penalty on β
- * Gradient descent w/ early stopping

RIDGE REGRESSION

$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left[\|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right]$

$\hat{\beta} = \underset{\beta}{\operatorname{argmin}}$

$$Y = X\beta + \epsilon$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left[\underbrace{\|Y - X\beta\|_2^2}_{\text{ERROR or LOSS}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{PENALTY}} \right]$$

RIDGE REGRESSION

$$\hat{\beta} = (X^{\top} X + \lambda I)^{-1} X^{\top} Y$$

$$\hat{\beta} = X_{ridge}^{+} Y$$

$$\hat{\beta} = (X^{\top} X + \lambda I)^{-1} X^{\top} Y$$

RIDGE REGRESSION

* Efficient solution with SVD

$$\hat{\beta} = (X^{\top} X + \lambda I)^{-1} X^{\top} Y$$

$$\text{(SVD)} \quad X = U S V^{\top} \quad D = \frac{S}{S^2 + \lambda^2}$$

$$\hat{\beta} = V D U^{\top} Y$$

$$\hat{\beta} = (X^{\top} X + \lambda I)^{-1} X^{\top} Y$$

$$D = \frac{S}{S^2 + \lambda^2}$$
$$\hat{\beta} = V D U^{\top} Y$$

RIDGE REGRESSION

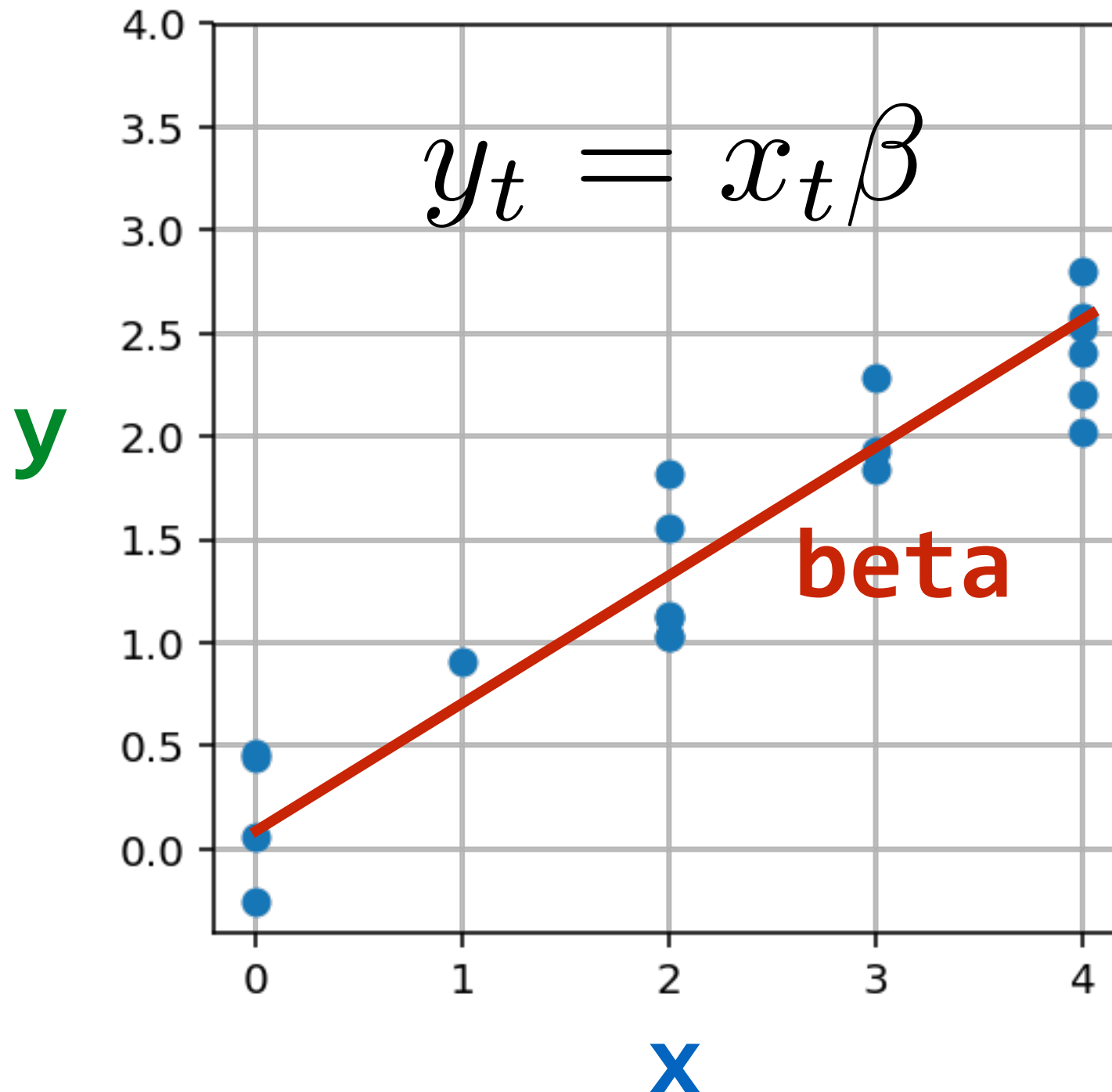
- * How to choose λ ?
- * GCV - Generalized Cross Validation 🙄
- * Block-wise cross-validation 👍

RIDGE REGRESSION

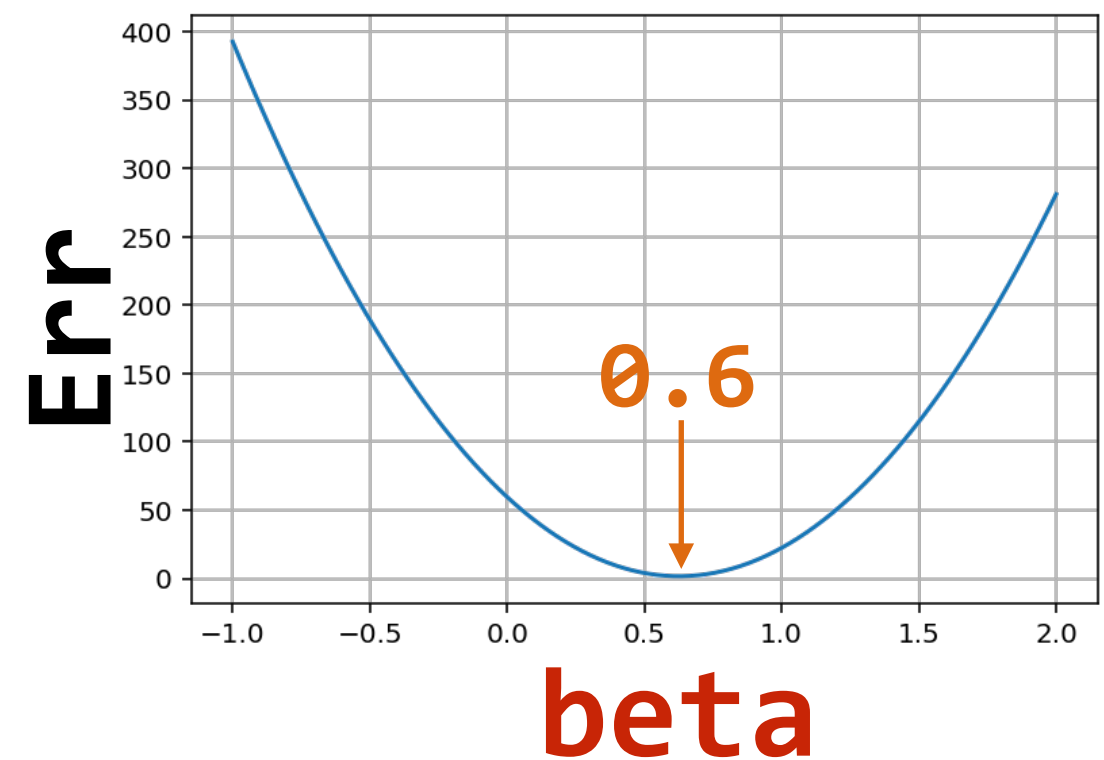
- * Good implementation: `scikit-learn`
- * Awesome implementation:
<http://github.com/alexhuth/ridge>

1D EXAMPLE

$$\text{Err}(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2$$



$$\text{Err}(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2$$

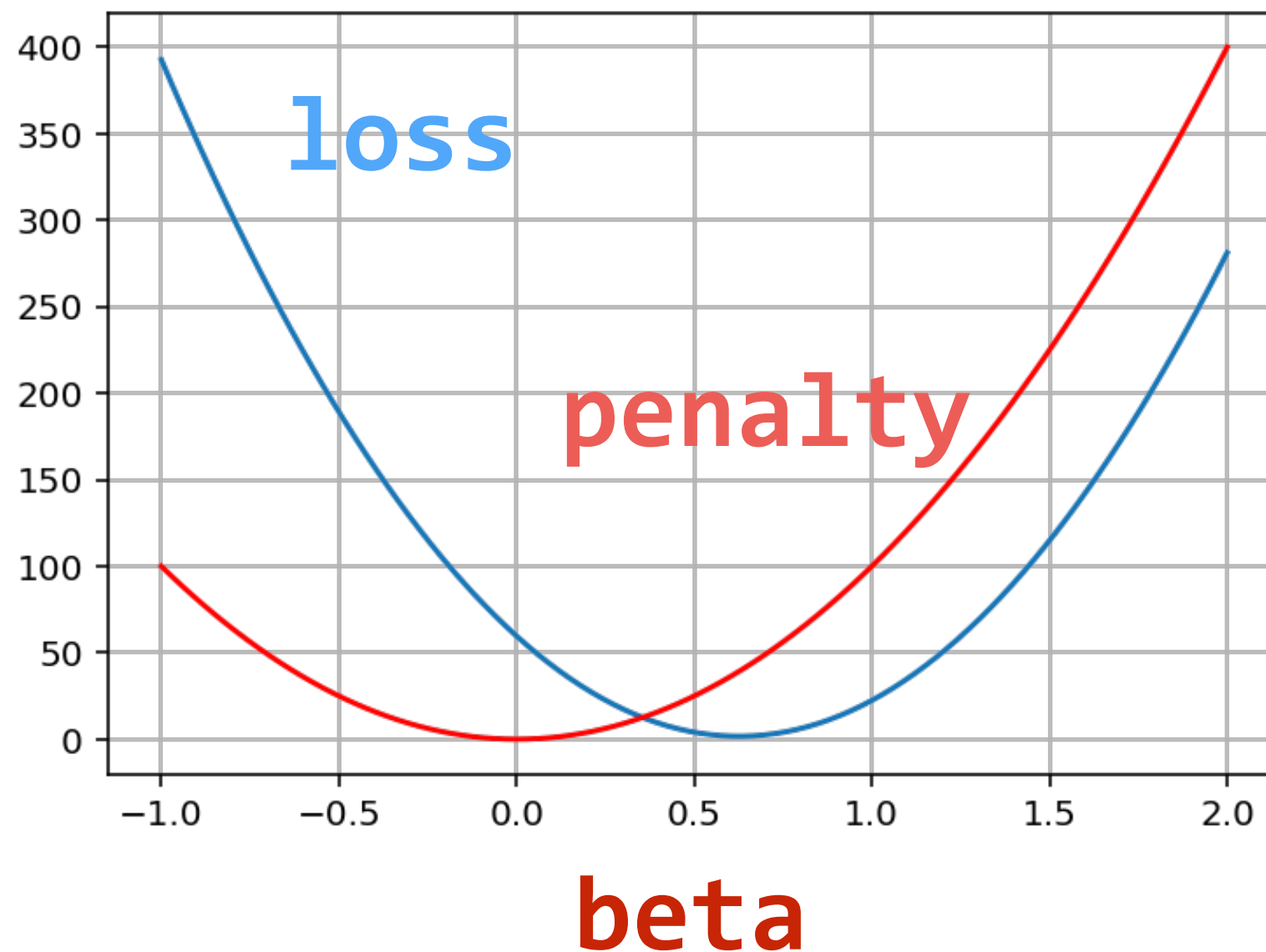


$$Err(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2 + \lambda \beta^2$$

1D EXAMPLE

L2 Regularization:
(as penalty)

$$Err(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2 + \lambda \beta^2$$

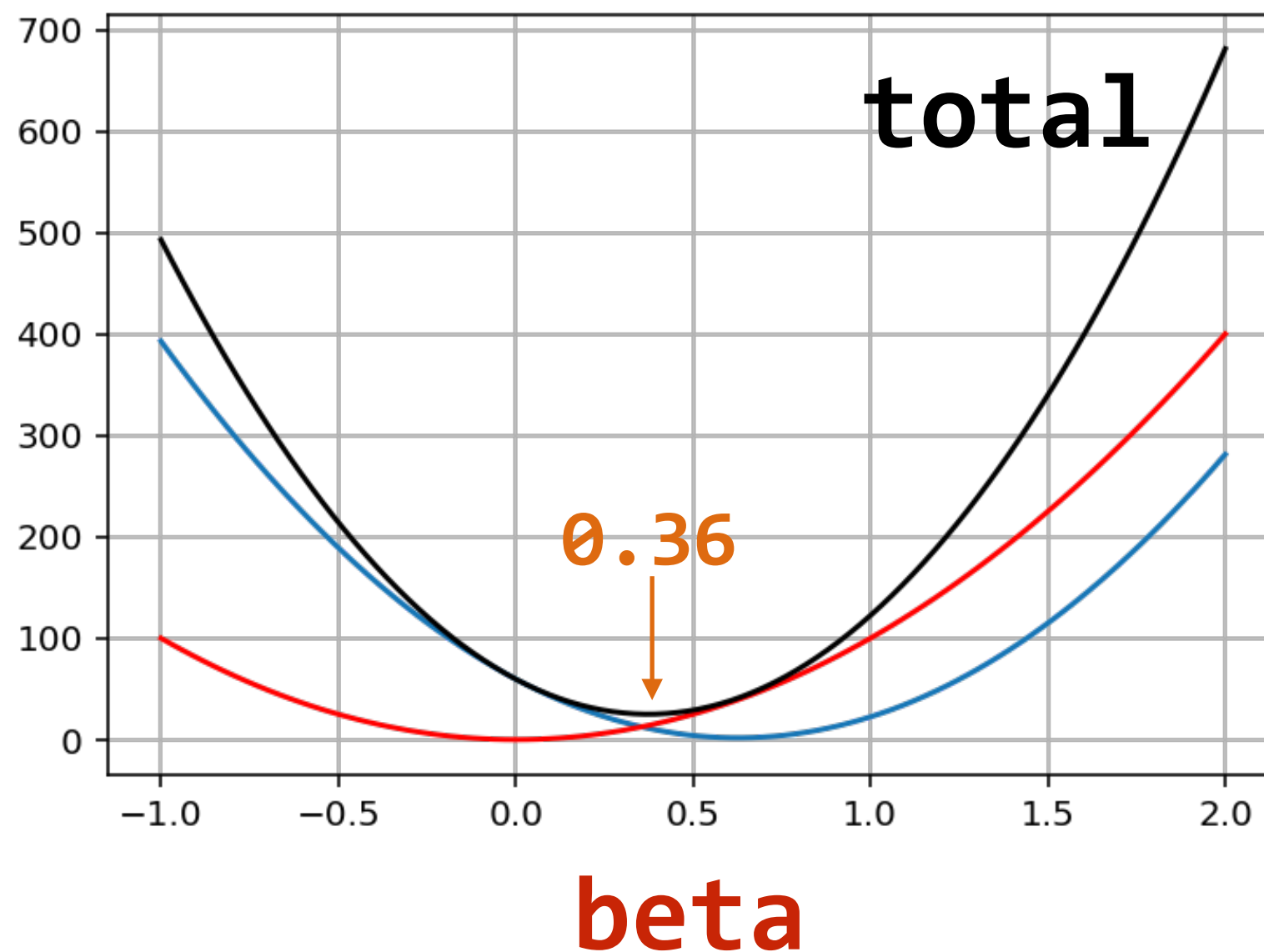


$$\text{Err}(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2 + \lambda \beta^2$$

1D EXAMPLE

L2 Regularization:
(as penalty)

$$\text{Err}(\beta) = \sum_{t=1}^T (y_t - x_t \beta)^2 + \lambda \beta^2$$



COVARIANCE

- * A measure of the “joint variability” of two variables

COVARIANCE

$$\text{cov}(\vec{x}, \vec{y}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

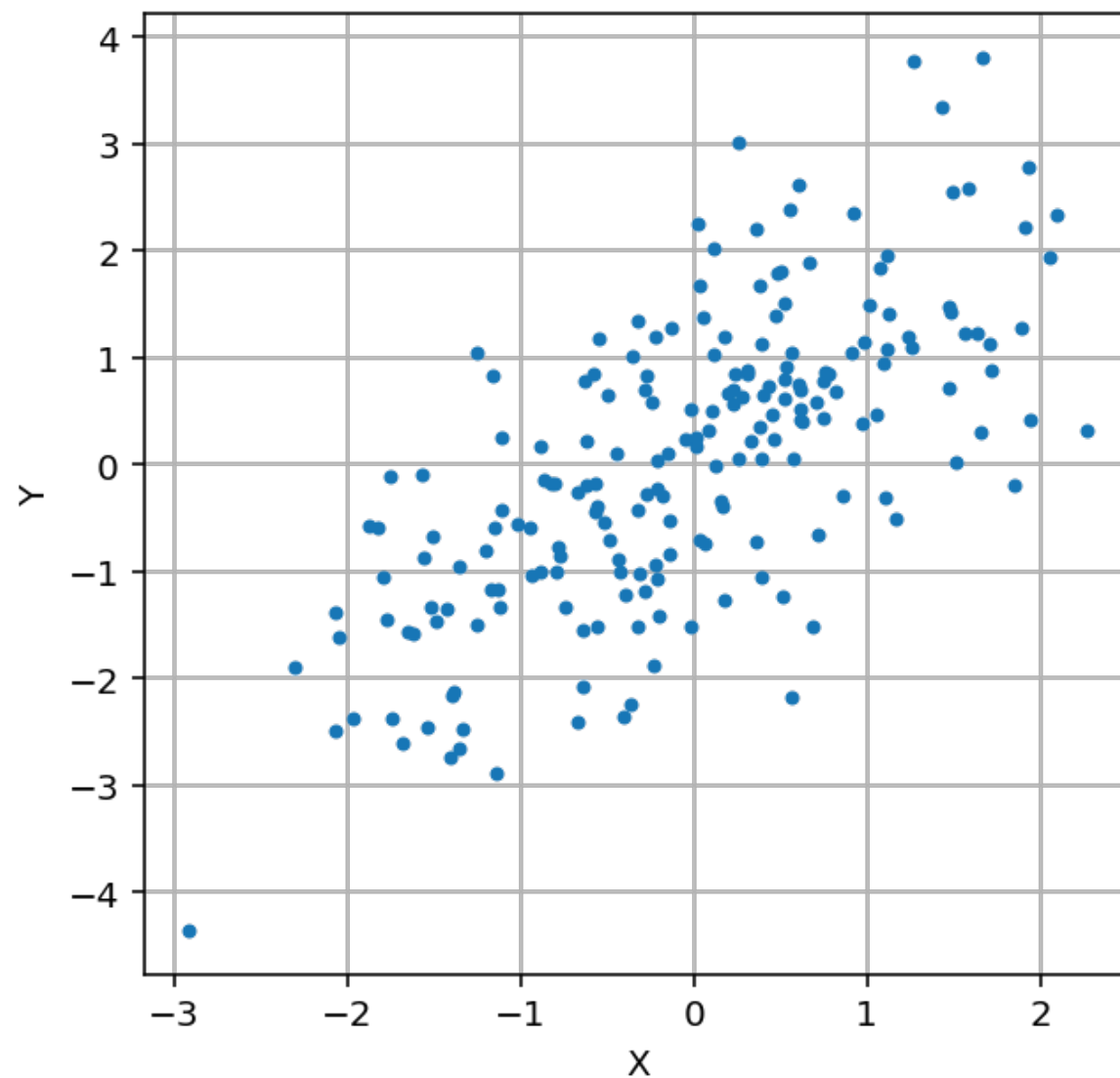
$$\text{cov}(X, X) = E[(X - E[X])(X - E[X])] = E[(X - E[X])^2] = \text{var}(X)$$

$$\text{cov}(\vec{x}, \vec{y}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$\text{cov}(X, X) = E[(X - E[X])(X - E[X])] = E[(X - E[X])^2] = \text{var}(X)$$

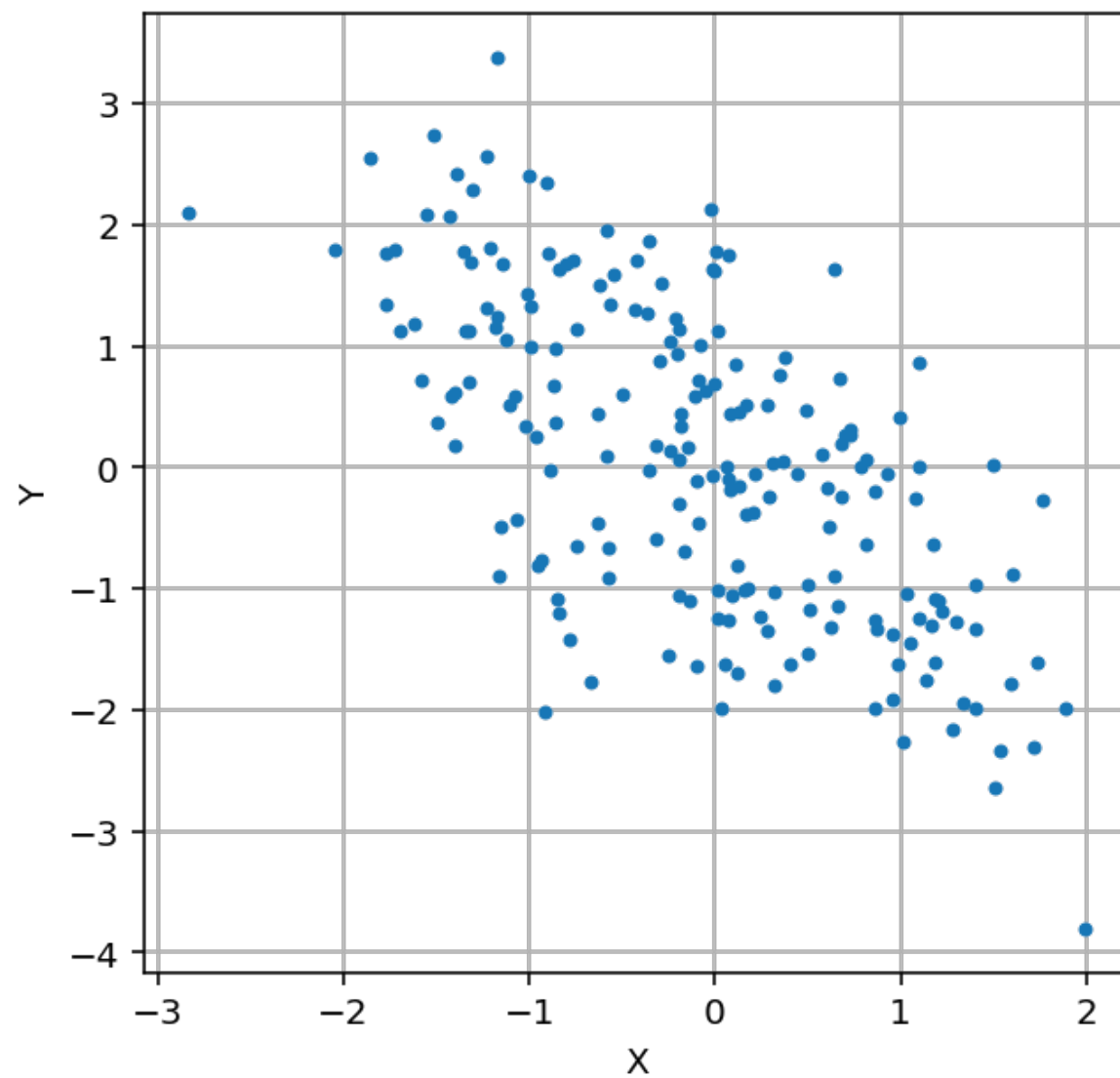
COVARIANCE



$$\text{cov}(x, y) > 0?$$

$$\text{cov}(x, y) < 0?$$

COVARIANCE

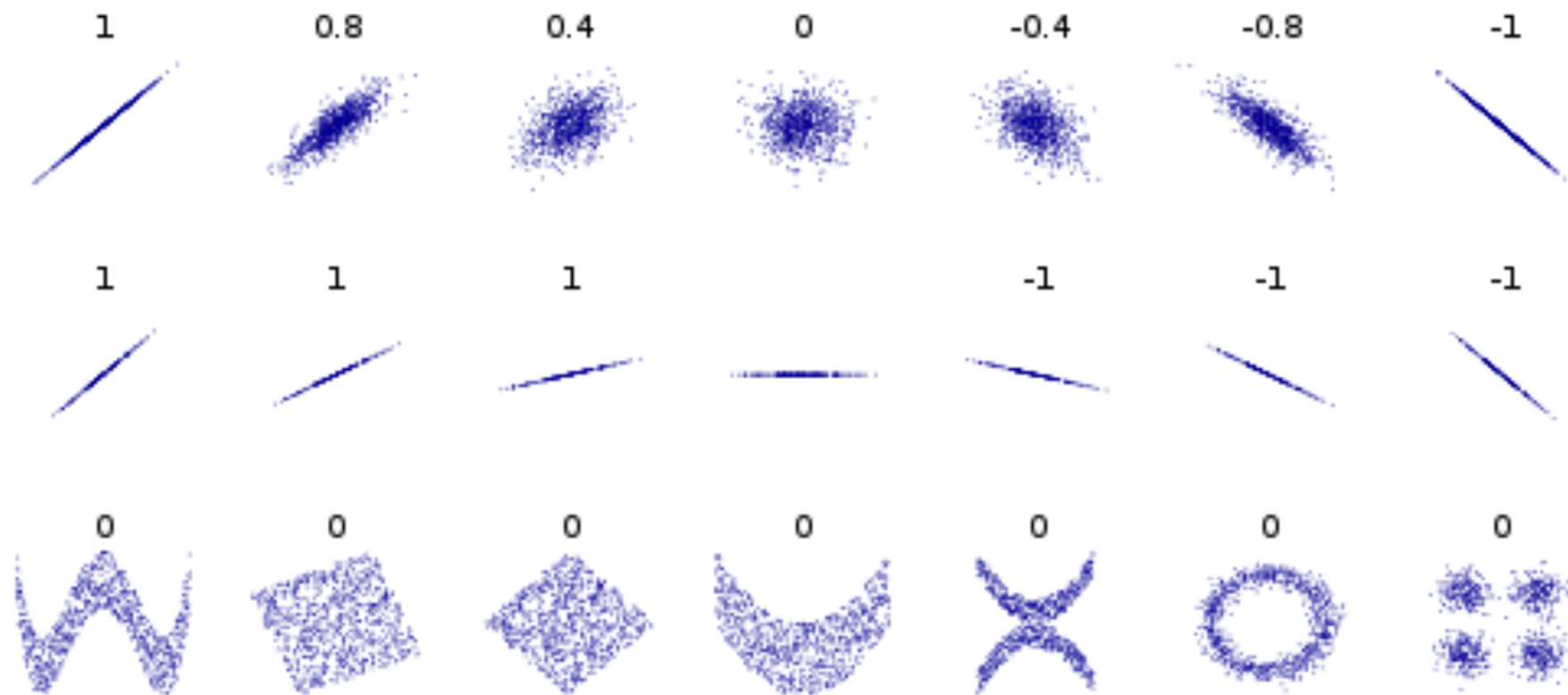


$$\text{cov}(x, y) > 0?$$

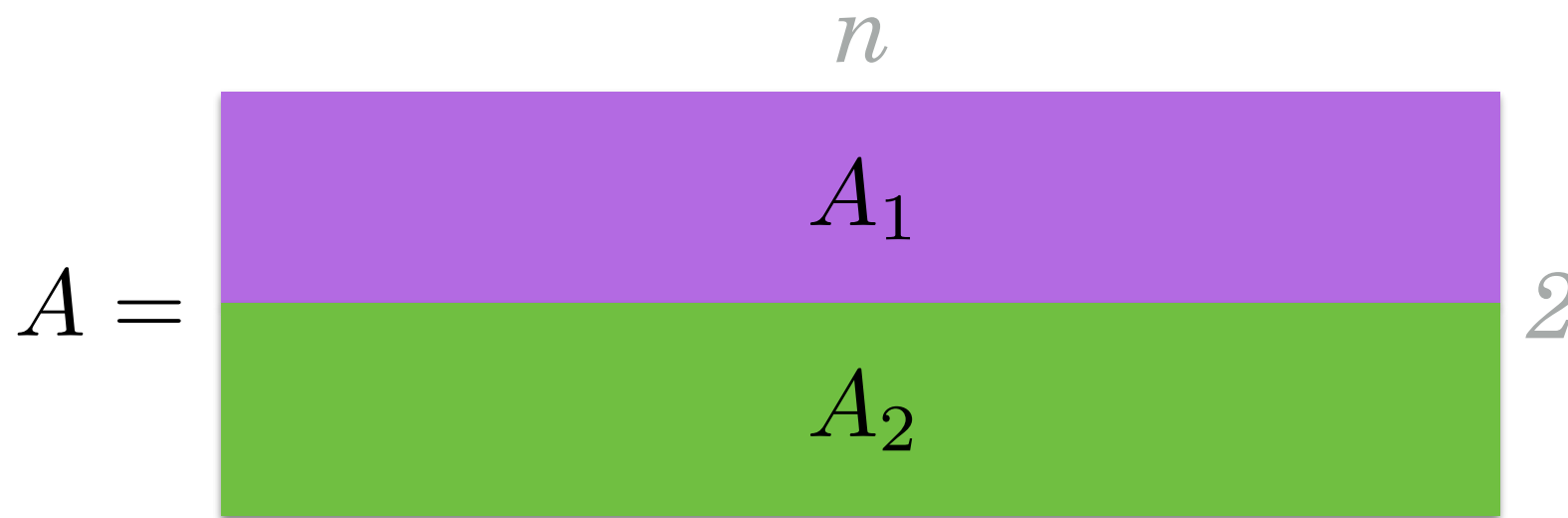
$$\text{cov}(x, y) < 0?$$

COVARIANCE

(Well, correlation...)



COVARIANCE MATRIX



```
\mbox{cov}(A) = \begin{bmatrix}
\mbox{var}(A_1) & \mbox{cov}(A_1, A_2) \\
\mbox{cov}(A_1, A_2) & \mbox{var}(A_2)
\end{bmatrix}
```

```
\mbox{cov}(A) = \left(\frac{1}{n}\right) A A^T
```

$$\text{cov}(A) = \begin{bmatrix} \text{var}(A_1) & \text{cov}(A_1, A_2) \\ \text{cov}(A_1, A_2) & \text{var}(A_2) \end{bmatrix}$$

(assuming A is mean 0)

$$\text{cov}(A) = \left(\frac{1}{n}\right) A A^T$$

TIKHONOV REGRESSION

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left[\|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right]$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}}$$

$$Y = X\beta + \epsilon$$

* RIDGE REGRESSION

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left[\underbrace{\|Y - X\beta\|_2^2}_{\text{ERROR or LOSS}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{PENALTY}} \right]$$

* TIKHONOV REGRESSION

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left[\|Y - X\beta\|_2^2 + \lambda \| \underset{\substack{\uparrow \\ \text{PENALTY} \\ \text{MATRIX}}}{C} \beta \|_2^2 \right]$$

TIKHONOV REGRESSION

- * **RIDGE REGRESSION** is a special case of **TIKHONOV REGRESSION**
- * **TIKHONOV REGRESSION** puts a **ZERO-MEAN MULTIVARIATE NORMAL PRIOR** on the weights
- * in **RIDGE REGRESSION** the covariance matrix of the prior has a constant diagonal
 - * i.e. the prior is a **SPHERE**
- * in **TIKHONOV REGRESSION** the covariance matrix can be ***ANYTHING***

TIKHONOV REGRESSION

- * the multivariate normal prior given by
TIKHONOV REGRESSION

$$\beta \sim N(0, \sigma^2 (C^T C)^{-1})$$

$\beta \sim N(0, \sigma^2 \Lambda^{-1}), \Lambda = C^T C$

TIKHONOV REGRESSION

```
\begin{eqnarray*}
A &=& XC^{-1} \\
\hat{\beta}_A &=& \underset{\beta}{\operatorname{argmin}} \left[ ||Y - A\beta||_2^2 + \lambda ||\beta||_2^2 \right]
\end{eqnarray*}
```

- * any **TIKHONOV** problem can be converted into a **RIDGE** problem

$$A = XC^{-1} \leftarrow 1. \text{ CHANGE OF BASIS}$$

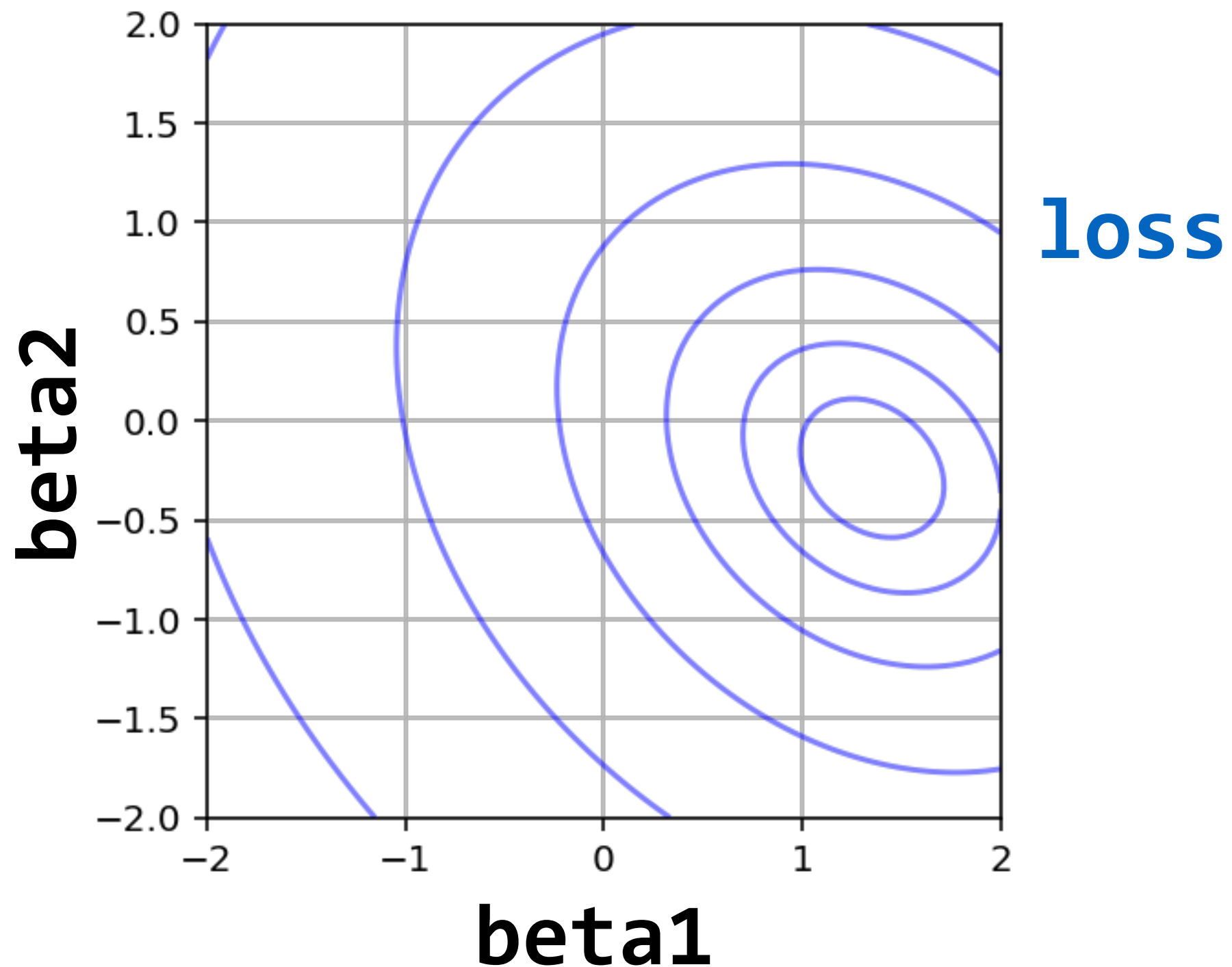
$$\hat{\beta}_A = \underset{\beta}{\operatorname{argmin}} \left[||Y - A\beta||_2^2 + \lambda ||\beta||_2^2 \right]$$

↑
2. RIDGE REGRESSION

$$\hat{\beta} = C^{-1} \hat{\beta}_A \leftarrow 3. \text{ CHANGE BASIS AGAIN}$$

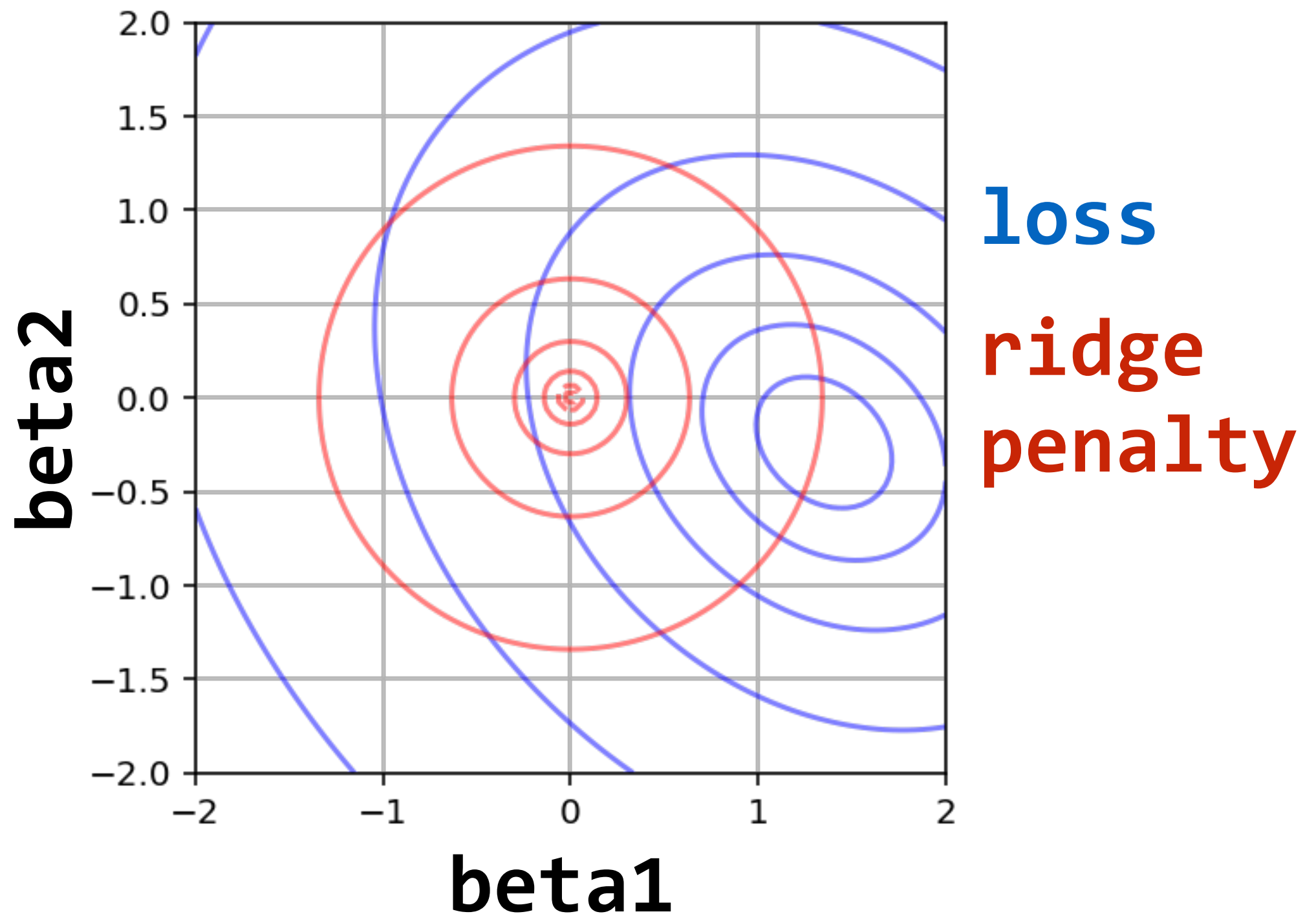
EXAMPLE

$$y_t = x_{1,t}\beta_1 + x_{2,t}\beta_2$$



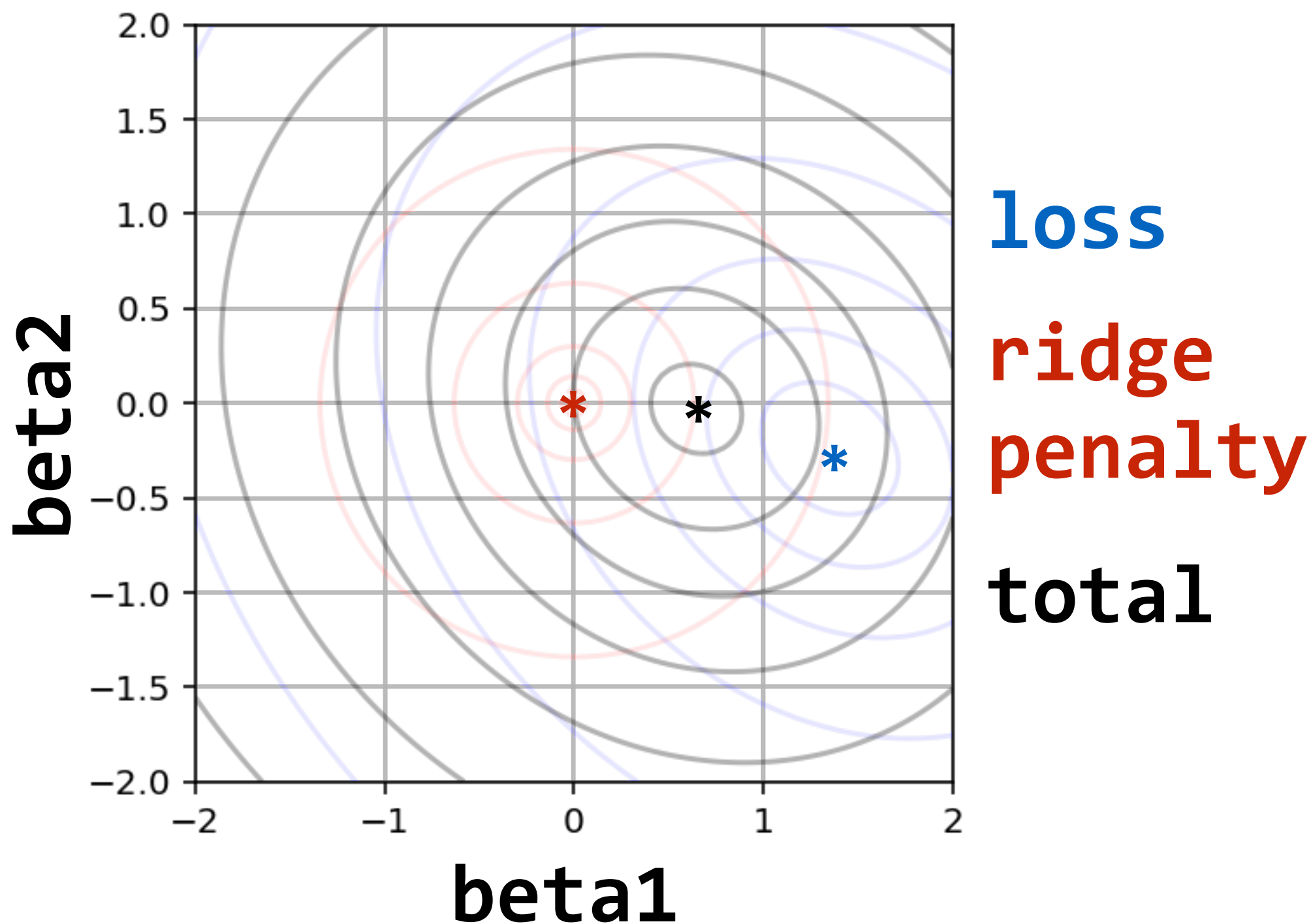
EXAMPLE

$$y_t = x_{1,t}\beta_1 + x_{2,t}\beta_2$$



EXAMPLE

$$y_t = x_{1,t}\beta_1 + x_{2,t}\beta_2$$

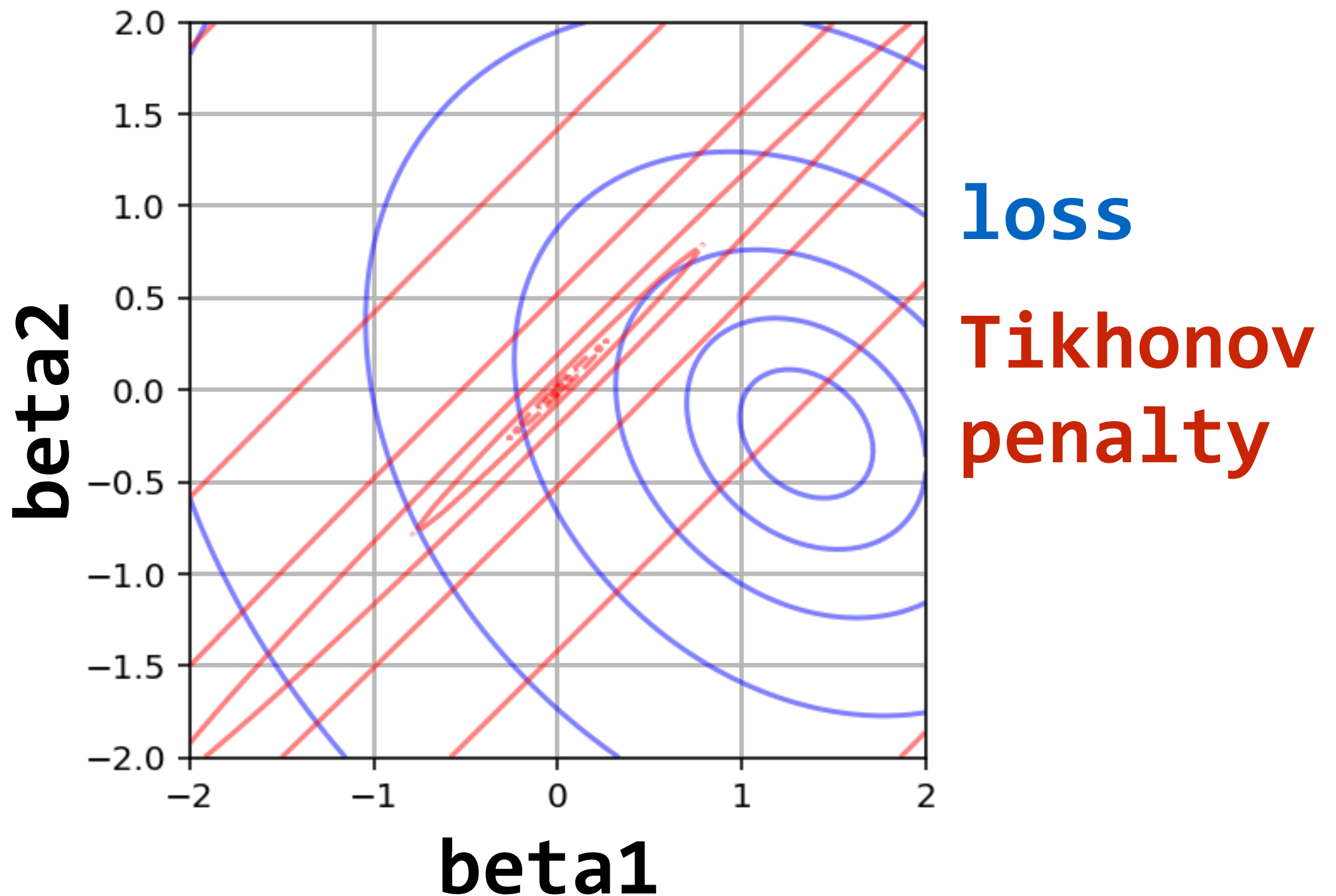


EXAMPLE

- * Suppose we strongly suspect that **beta1** and **beta2** should be similar

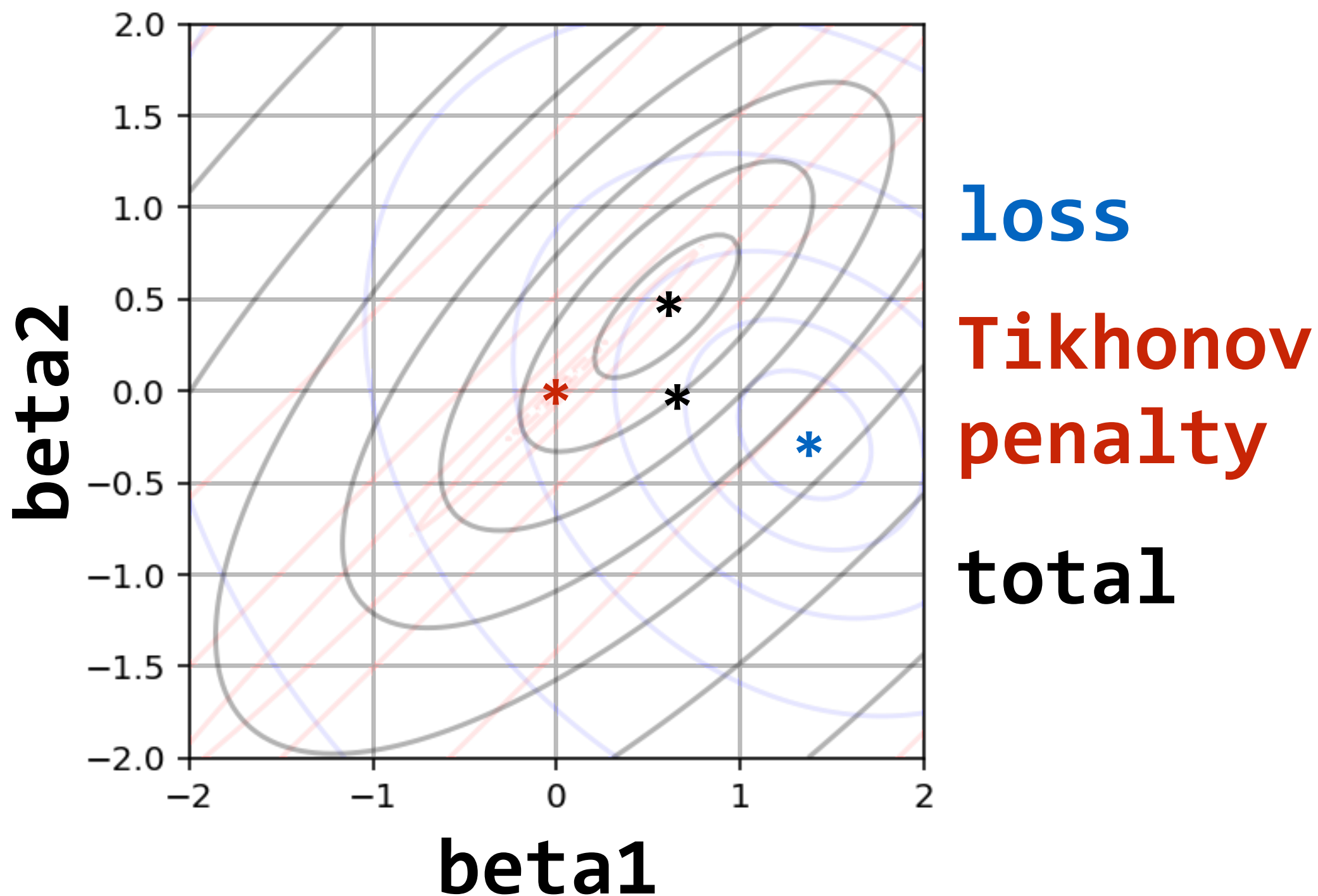
EXAMPLE

$$y_t = x_{1,t}\beta_1 + x_{2,t}\beta_2$$



EXAMPLE

$$y_t = x_{1,t}\beta_1 + x_{2,t}\beta_2$$



TIKHONOV REGRESSION

- * any **TIKHONOV** problem can be converted into a **RIDGE** problem by a **LINEAR TRANSFORMATION**
- * conversely, **ANY LINEAR TRANSFORMATION** of X followed by **RIDGE REGRESSION** is equivalent to some **TIKHONOV REGRESSION** problem

TIKHONOV REGRESSION

- * WORD EMBEDDING MODELS
- * think of stimulus matrix as **WORDS** over time projected onto **WORD EMBEDDING**

$$\begin{array}{c} \text{WORDS} \quad \text{EMBEDDING} \\ \downarrow \quad \downarrow \\ X = W E \\ \begin{array}{ccc} t \times n & t \times w & w \times n \\ \uparrow & \uparrow & \uparrow \\ \text{TIME} & \text{WORDS} & \text{EMBEDDING DIMENSIONS} \end{array} \end{array}$$

$$\begin{array}{l} \text{\texttt{\textbackslashunderset\{t\}\times n\}\{X\}} = \\ \text{\texttt{\textbackslashunderset\{t\}\times w\}\{W\}} \\ \text{\texttt{\textbackslashunderset\{w\}\times n\}\{E\}} \end{array}$$

TIKHONOV REGRESSION

- * this is equivalent to **TIKHONOV REGRESSION** on the **WORDS** with a prior determined by the **WORD EMBEDDING**

$$\frac{1}{\sigma^2} \Sigma_{\beta} = (C^T C)^{-1} = E^T E$$

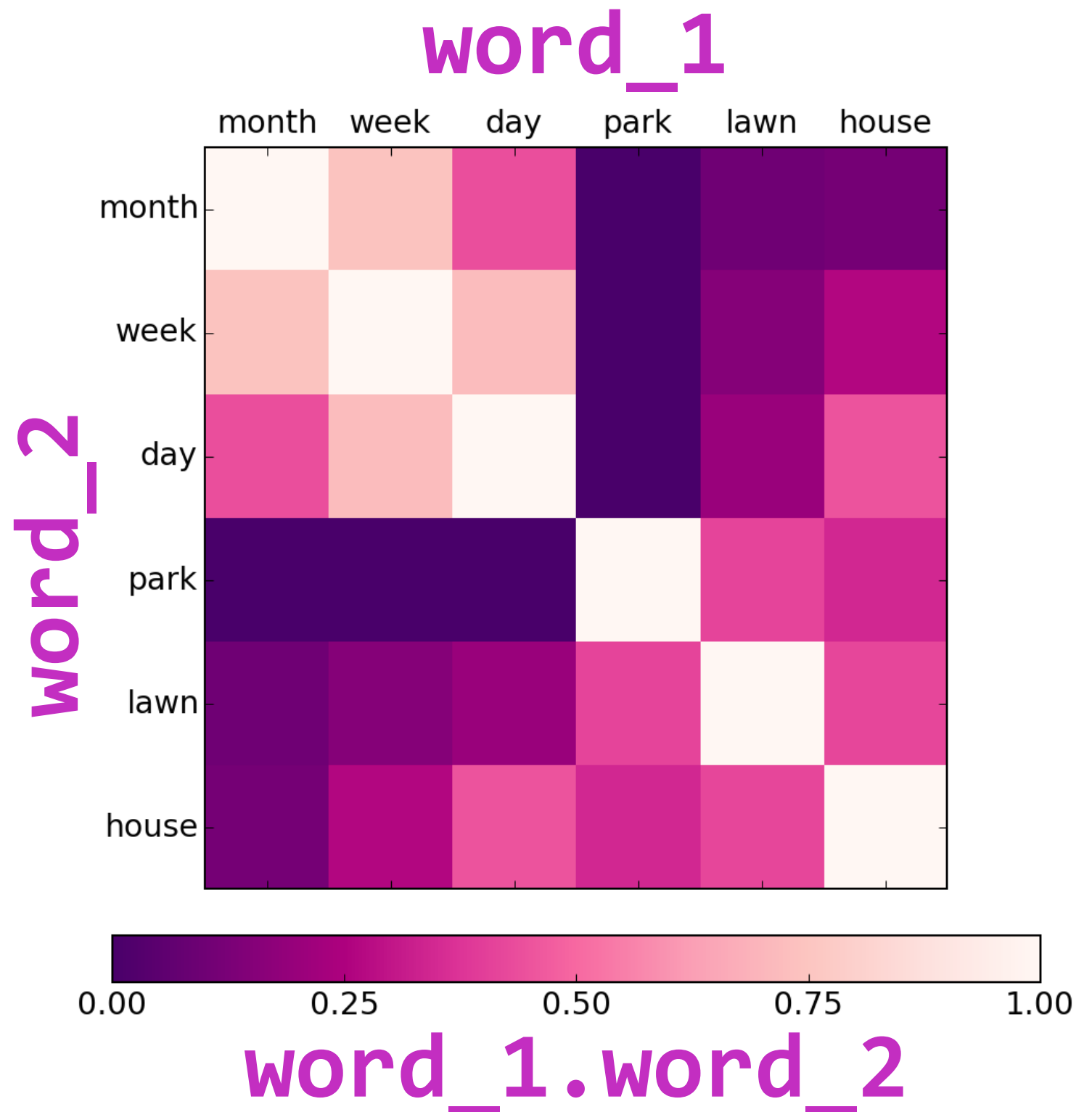
PRIOR
COVARIANCEINVERSE OF
PENALTY
INNER PRODUCTEMBEDDING
INNER PRODUCT

- * i.e. the prior covariance between two words' weights is equal to the dot product of their embedding vectors

$\frac{1}{\sigma^2} \Sigma_{\beta} = E^T E$

TIKHONOV REGRESSION

$E^T E =$
EMBEDDING
INNER PRODUCT,
english1000



TIKHONOV REGRESSION

$\underset{w \times v}{\text{}} \\ \underset{w \times n}{\text{}} \\ \underset{n \times v}{\text{}}$

- * to get **WEIGHTS ON WORDS** we just project onto the **EMBEDDING**

WEIGHTS IN
WORD SPACE

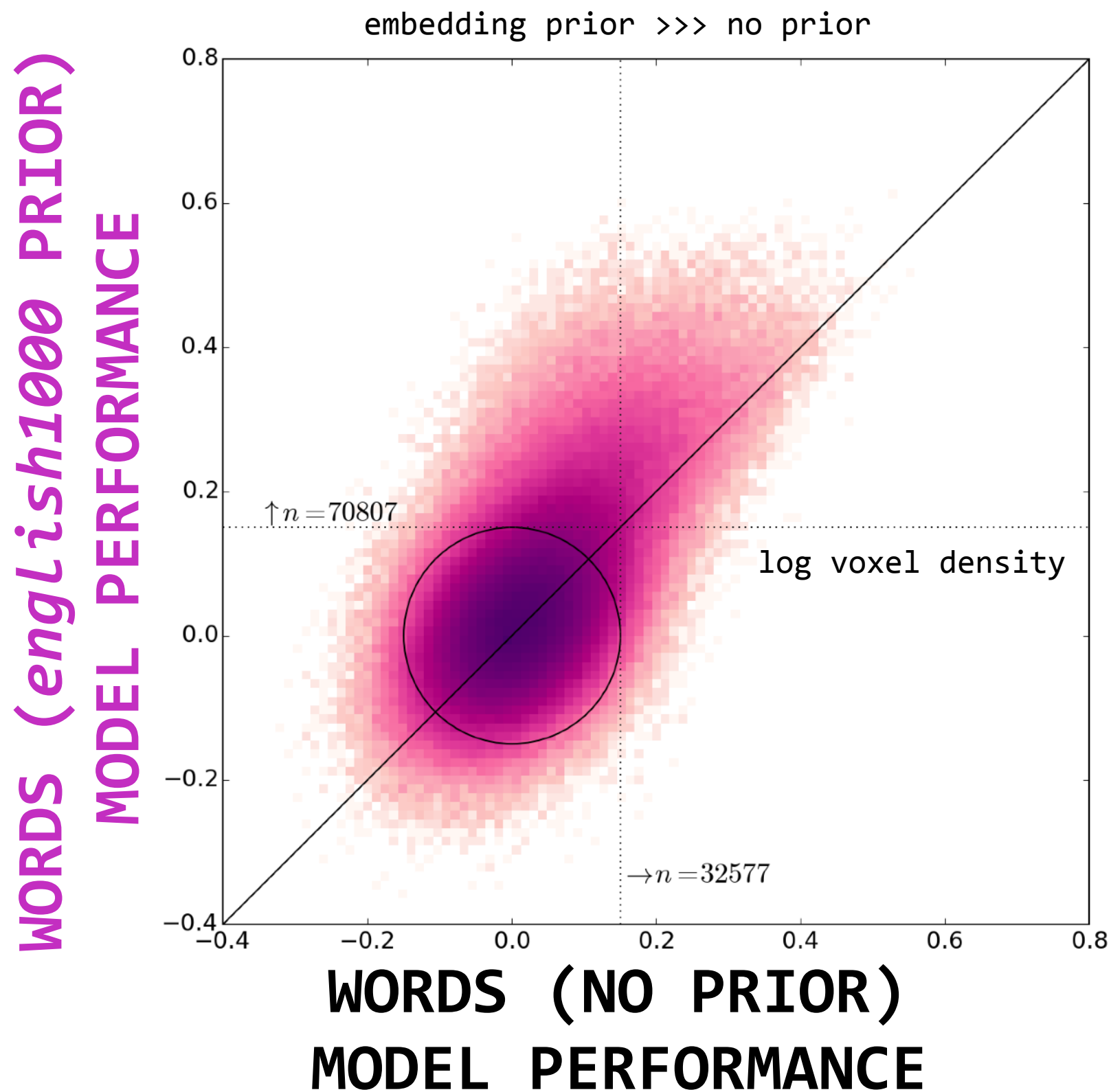
EMBEDDING

WEIGHTS IN
EMBEDDING SPACE

$$\underset{w \times v}{\hat{\beta}_W} = \underset{w \times n}{E} \underset{n \times v}{\hat{\beta}_X}$$

- * (this is equivalent to simulating responses to single words)

TIKHONOV REGRESSION



NEXT TIME

- * Data quality!