

Overfitting and regularization

Herman Kamper

<http://www.kamperh.com/>

Linear regression

Examples of overfitting

Herman Kamper

<http://www.kamperh.com/>

Overfitting example

Suppose we want to fit a regression model with scalar input using basis functions. Also suppose we have $N=10$ training items.

Our goal: $\underline{y} \approx \underline{\Phi} \underline{w}$

If we use 2 basis functions, the shapes will be:

$$\begin{matrix} 10 \times 1 & & 10 \times 2 & 2 \times 1 \\ \underline{y} & \approx & \underline{\Phi} & \underline{w} \end{matrix}$$

If instead we have 10 basis functions

we would have:

$$\begin{matrix} 10 \times 1 & & 10 \times 10 & 10 \times 1 \\ \underline{y} & \approx & \underline{\Phi} & \underline{w} \end{matrix}$$

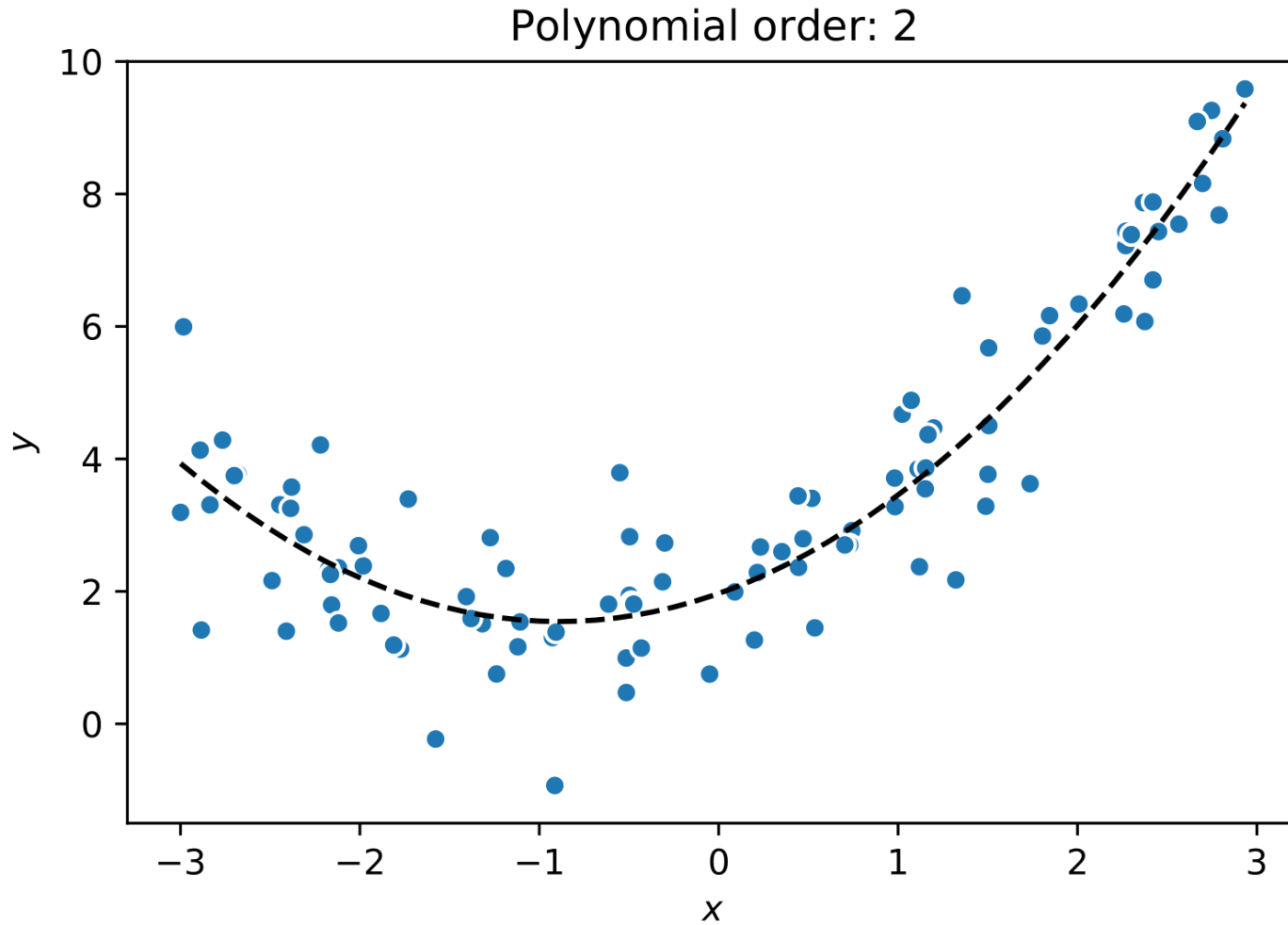
But this is solvable exactly! 10 equations in 10 unknowns (the 10 weights). So we can solve exactly: $\underline{w} = \underline{\Phi}^{-1} \underline{y}$

Questions

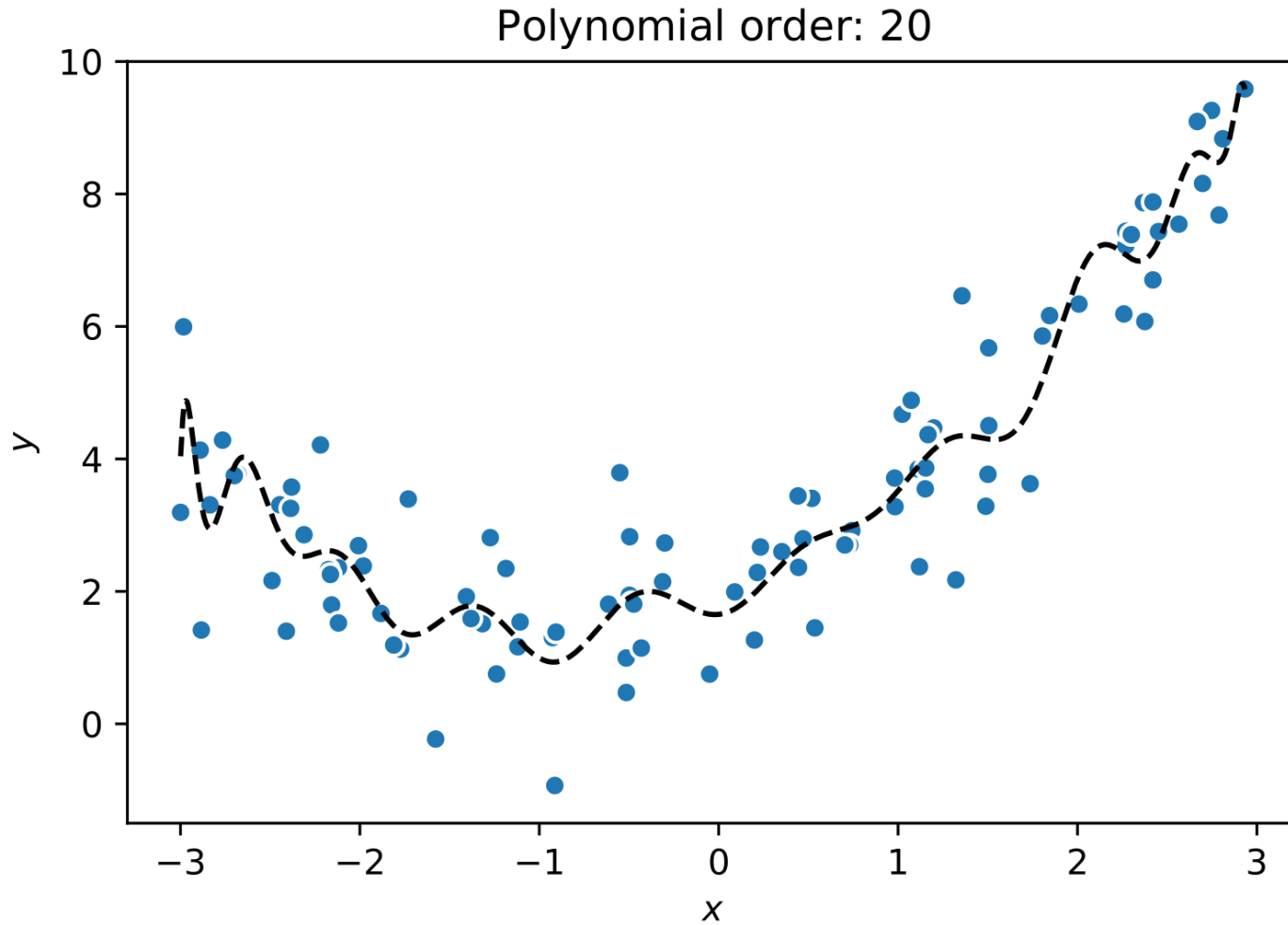
- What would the value of the loss J be?
- Would this be a good fit for making future predictions?

} If $\underline{\Phi}$ is invertible.

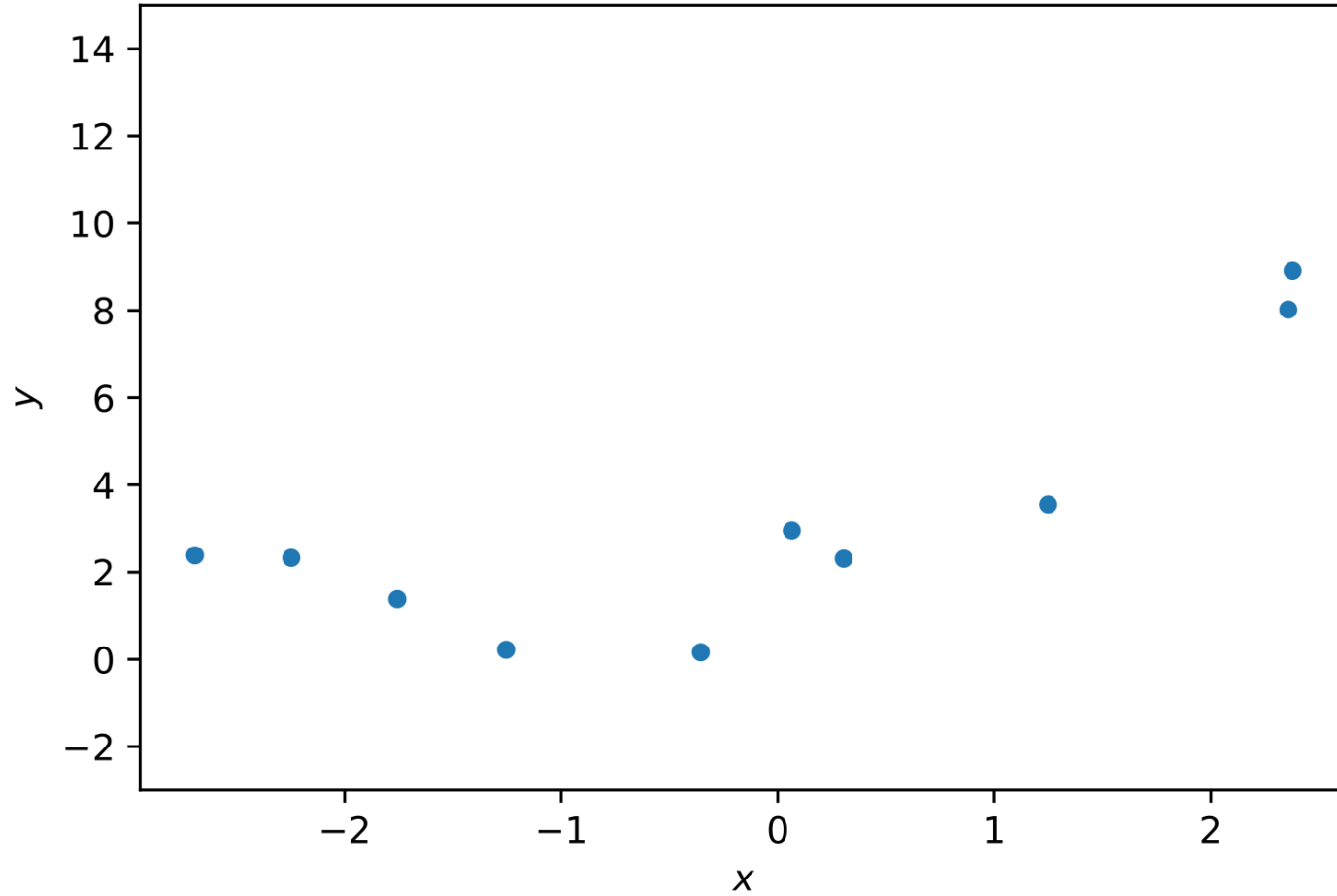
Polynomial regression



Polynomial regression



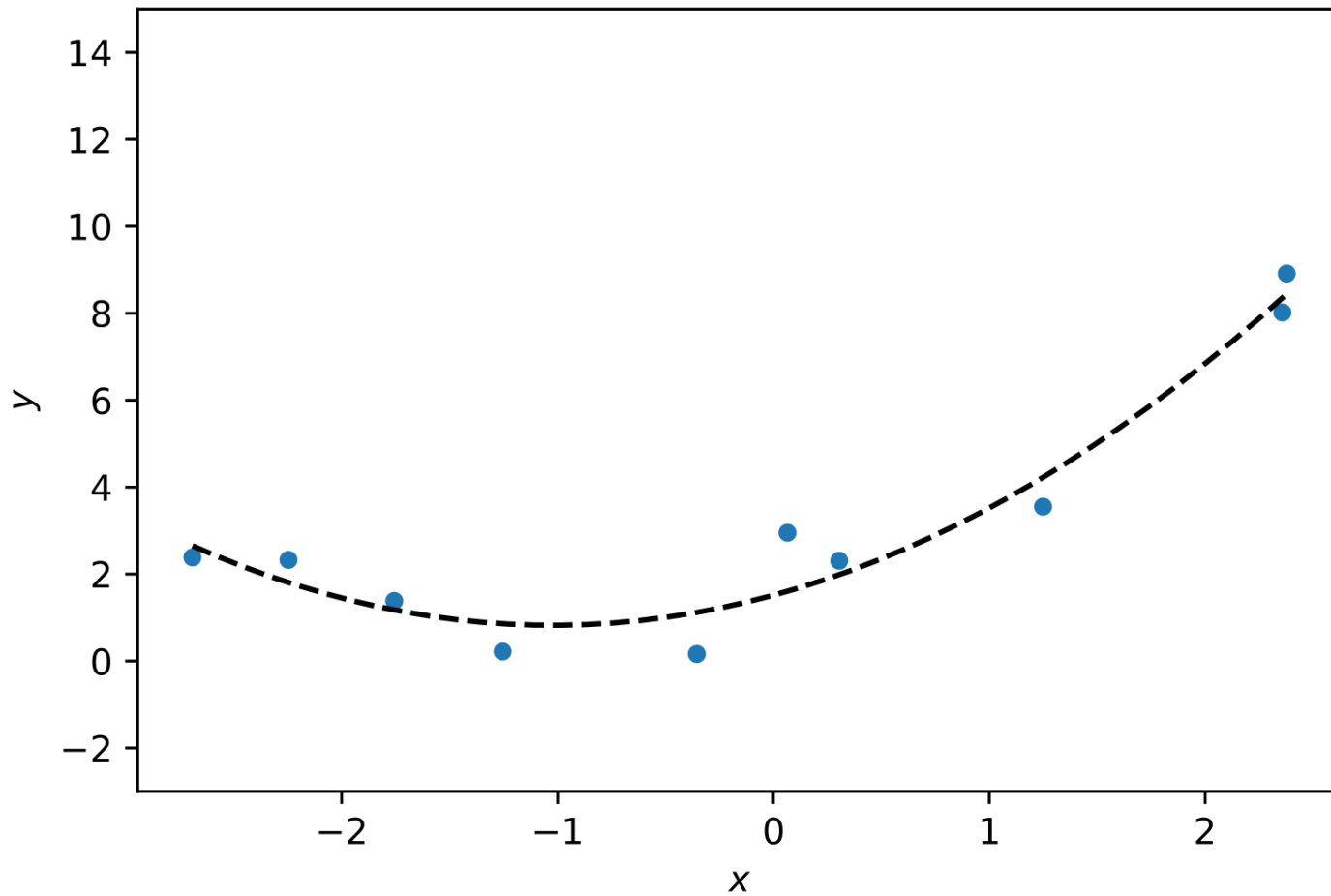
Polynomial regression



Polynomial regression

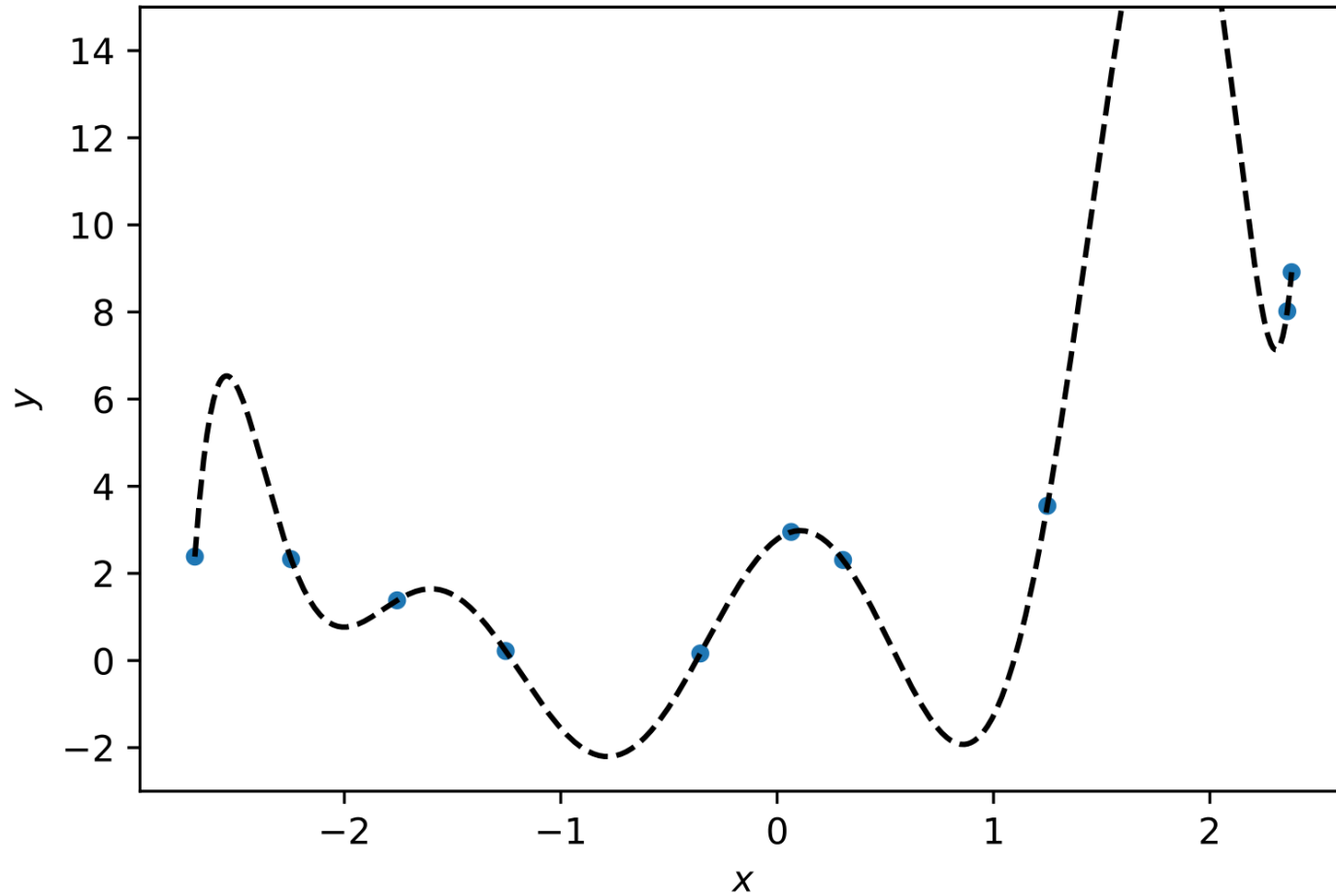
$$f(x; \hat{w}) =$$

$$f(x) = 1.15 + 1.35x + 0.66x^2$$



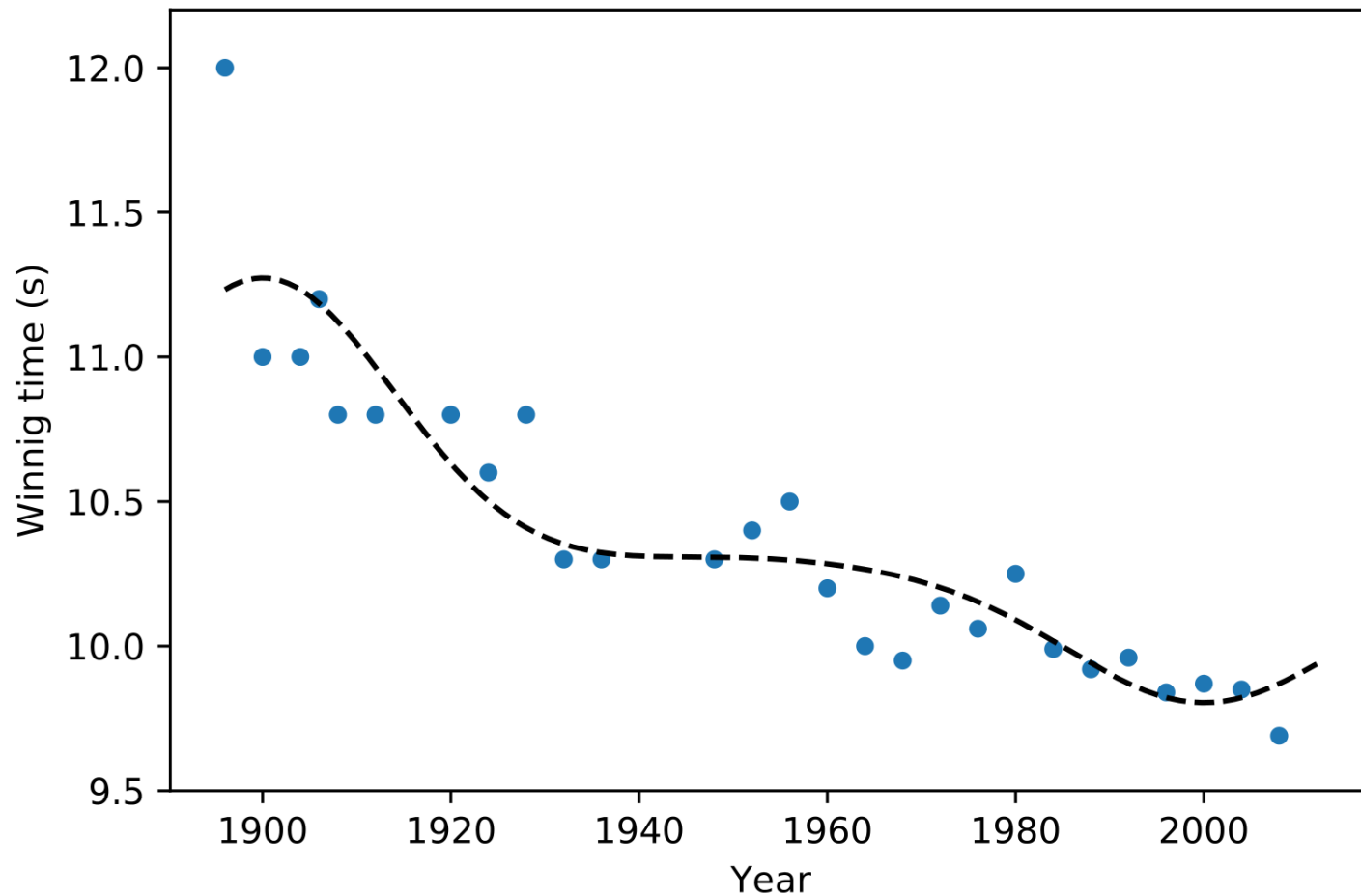
Polynomial regression

$$f(x) = 2.79 + 3.59x^2 - 15.61x^3 - 9.57x^4 + 15.11x^5 + 8.01x^6 - 4.03x^7 + \dots$$

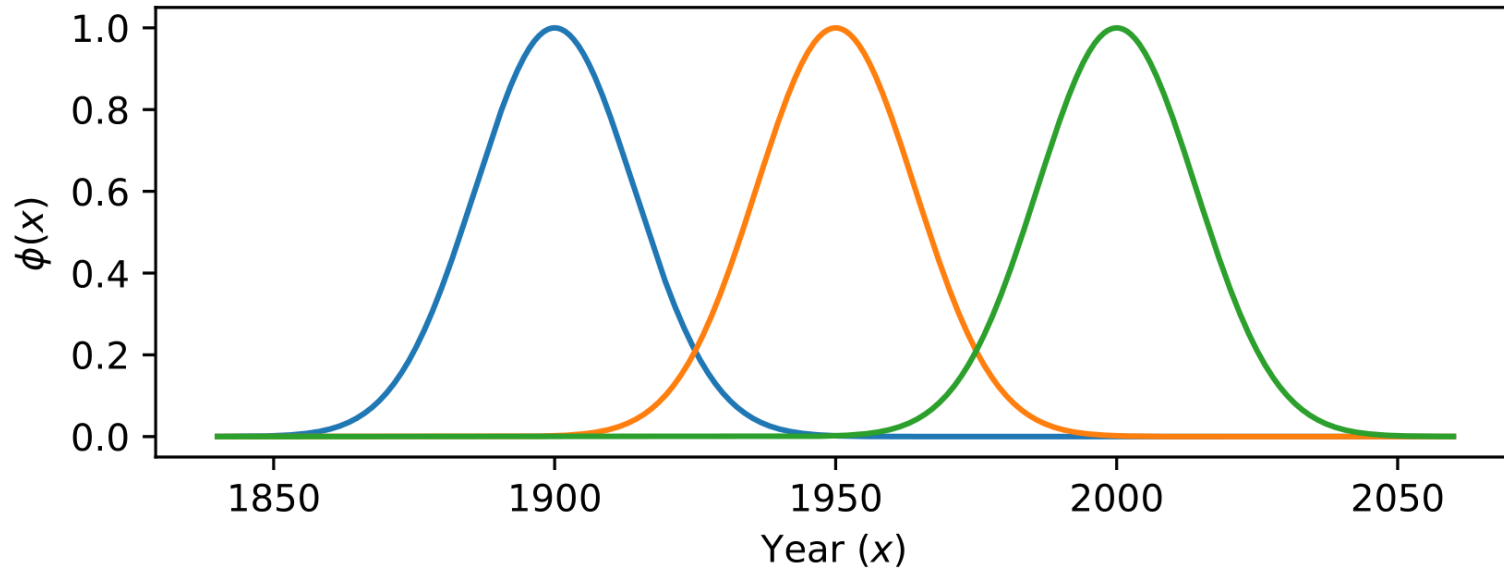


RBF with $c = [1900, 1950, 2000]$ and $h = 20$

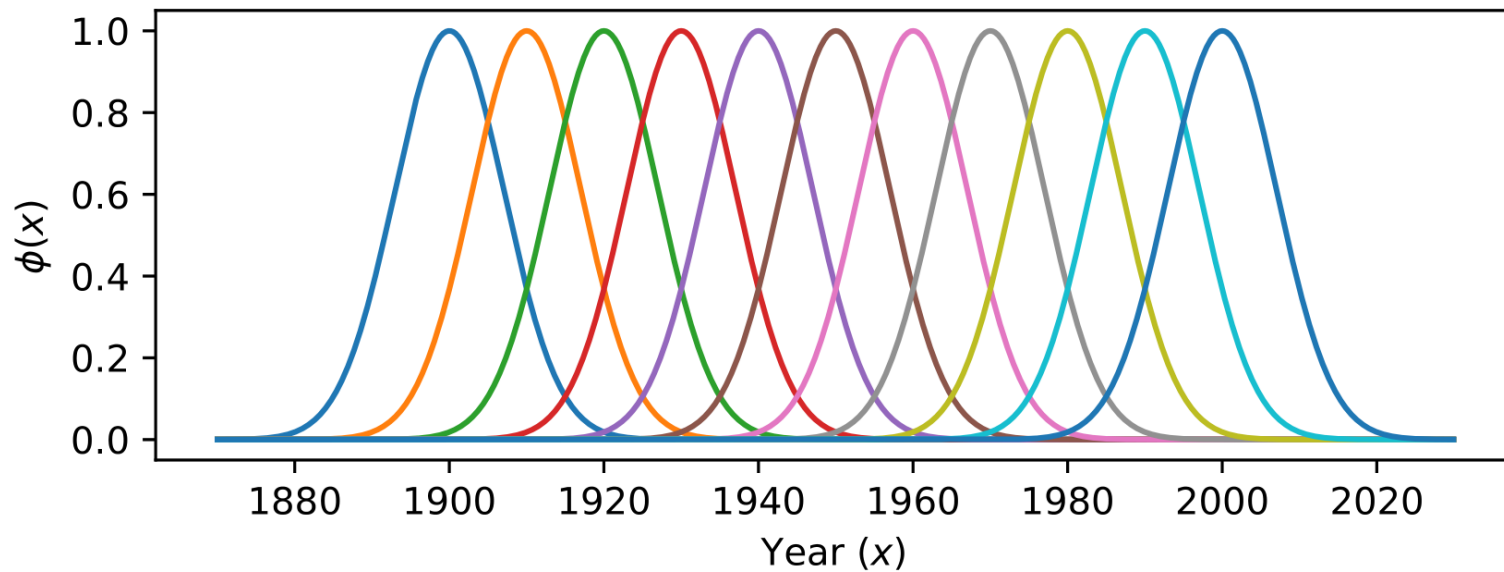
$$\sum_{k=1}^K w_k^2 = 1.25$$



RBF with $c = [1900, 1950, 2000]$ and $h = 20$

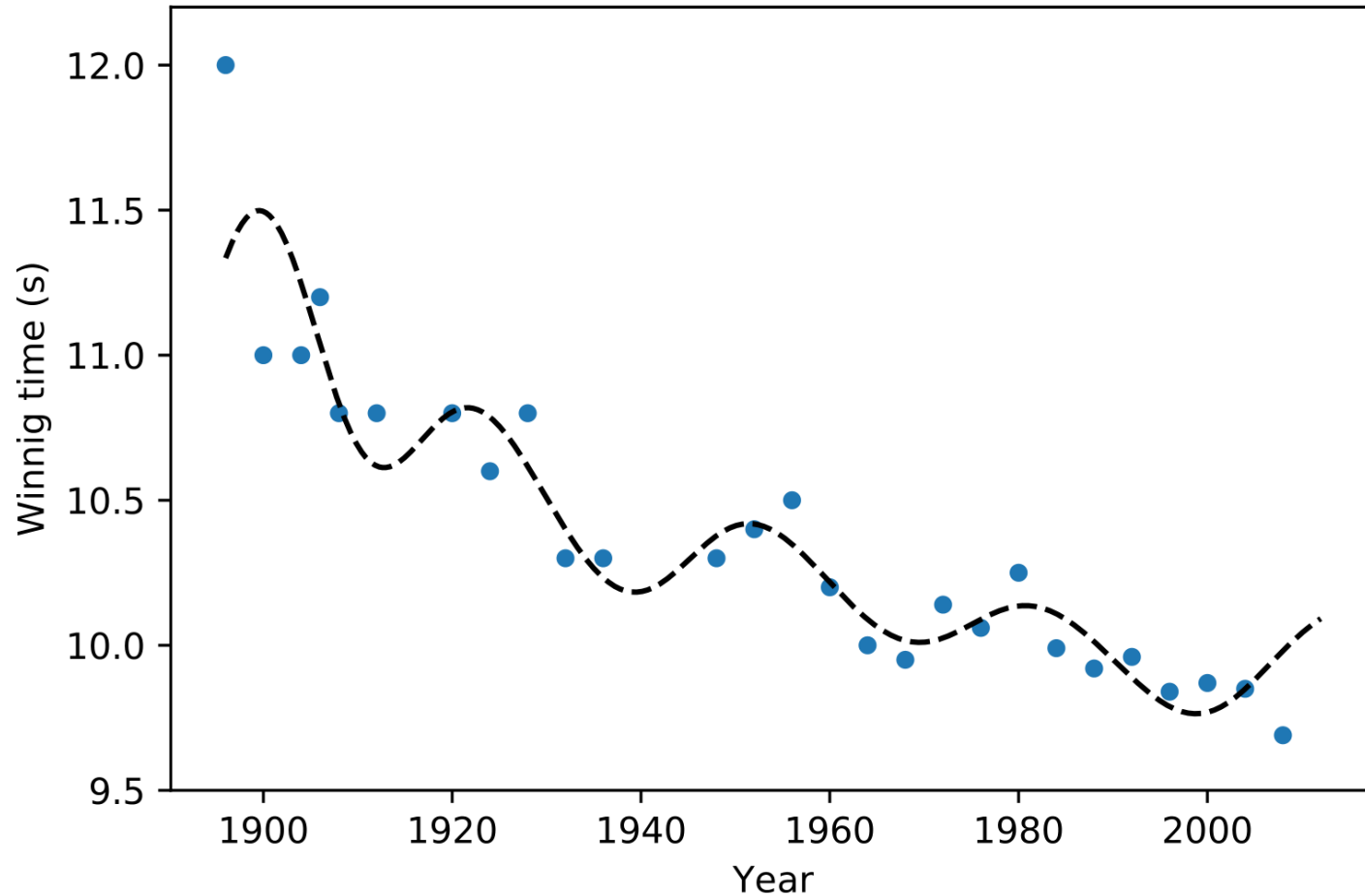


RBF with $c = [1900, 1910, \dots, 2000]$ and $h = 10$

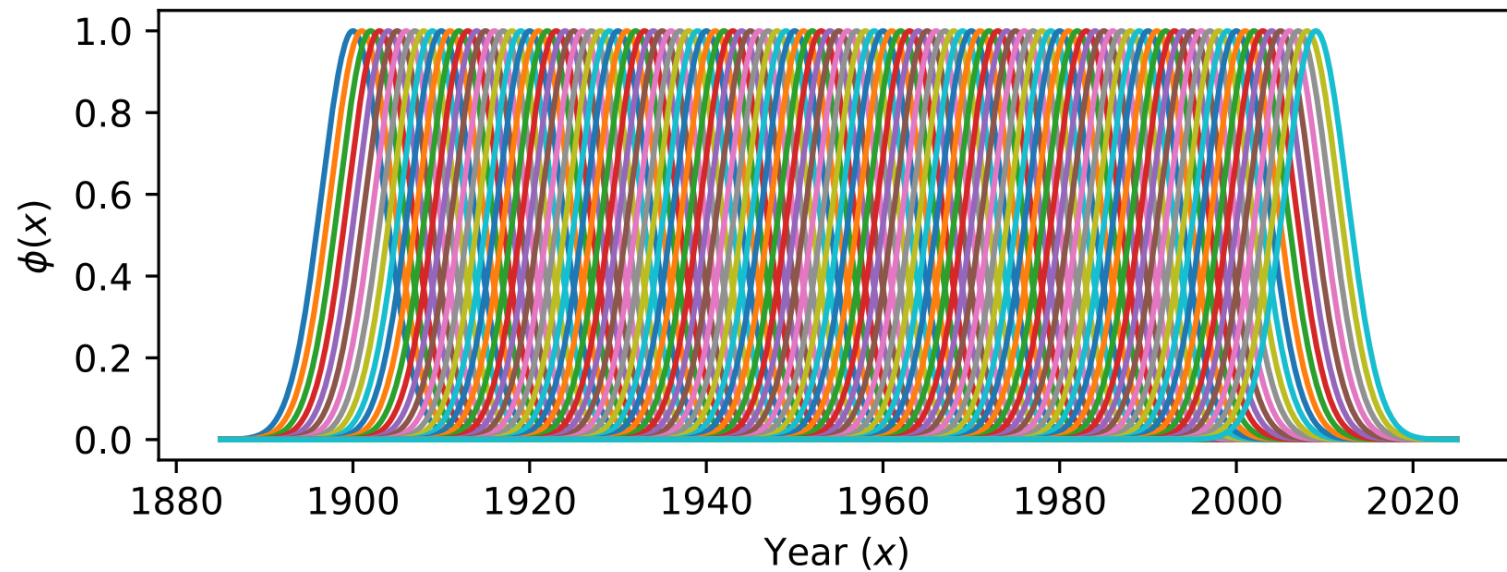


RBF with $c = [1900, 1910, \dots, 2000]$ and $h = 10$

$$\sum_{k=1}^K w_k^2 = 2.74$$

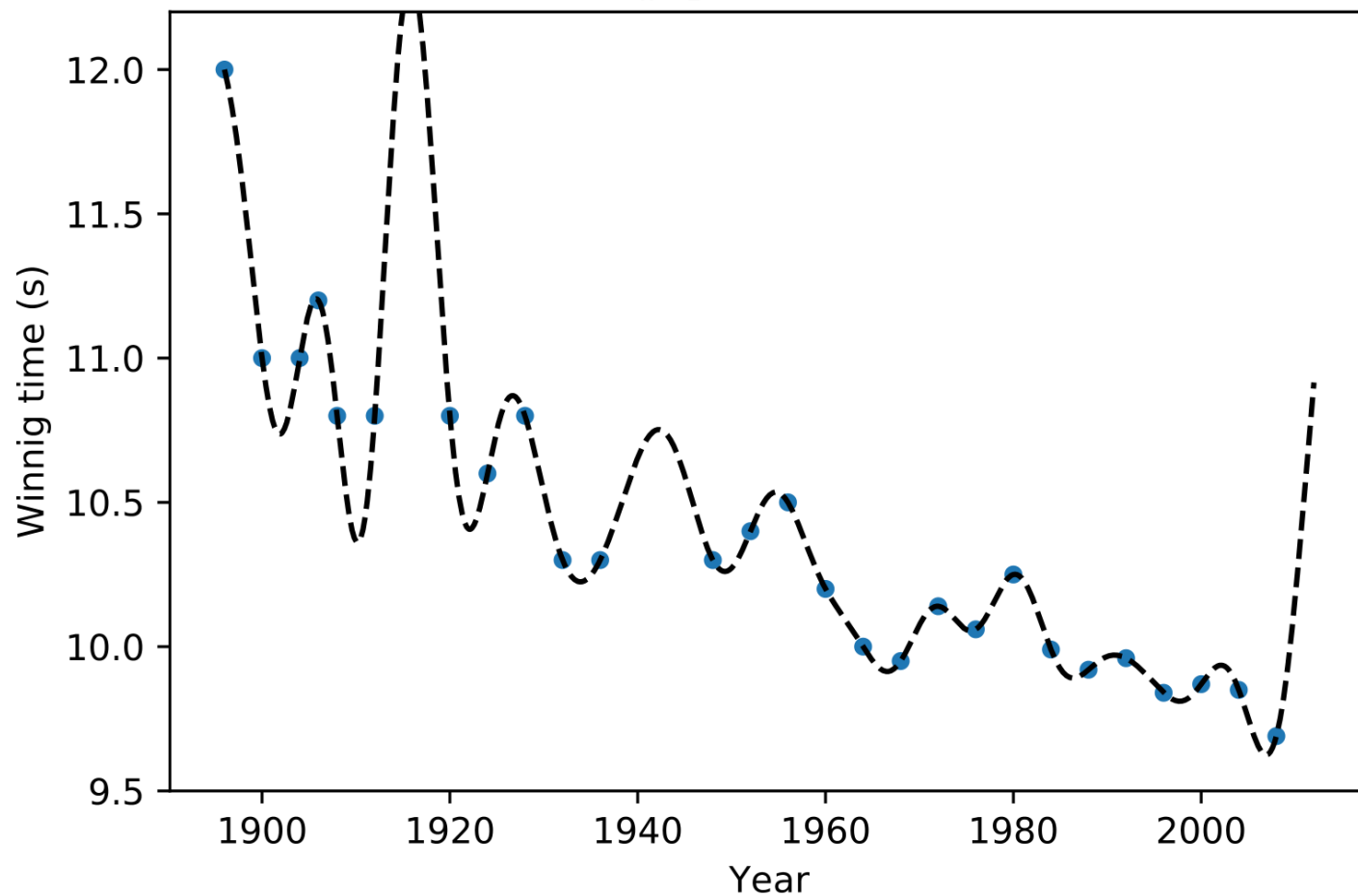


RBF with $c = [1900, 1901, \dots, 2000]$ and $h = 1$



RBF with $c = [1900, 1901, \dots, 2000]$ and $h = 1$

$$\sum_{k=1}^K w_k^2 = 20.79$$



Regularization

Combatting overfitting

Herman Kamper

<http://www.kamperh.com/>

Regularization:

Might want to fit higher-order models, but then it would be useful to be able to control their "complexity" in some way.

Idea:

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{arg\,min}} \left\{ \underset{\substack{\text{Standard loss (e.g. squared)}}{\mathcal{J}(\underline{w})}} + \text{penalty}(\underline{w}) \right\}$$

Penalty functions that constrains \underline{w} to be small are sometimes called "shrinkage" methods. We consider two penalty methods:

- Ridge (L_2) regularization
- Lasso (L_1) regularization

Ridge (L2) regularization: *Need to choose (see later)*

$$J_{\lambda}(\underline{w}) = \sum_{n=1}^N (y^{(n)} - f(\underline{x}^{(n)}; \underline{w}))^2 + \lambda \sum_{k=1}^K w_k^2 \quad \dots \textcircled{1}$$

We normally don't regularize w_0 . Why not? An easy hack is to zero-mean your data beforehand, i.e. the columns of \underline{X} (or $\underline{\Phi}$) normalized to have a mean of $\underline{0}$.

$$J_{\lambda}(\underline{w}) = \sum_{n=1}^N (y^{(n)} - f(\underline{x}^{(n)}; \underline{w}))^2 + \lambda \underline{w}^T \underline{w}$$

Can find closed-form solution exactly as before:

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

D-dimensional identity matrix

Note that anywhere we have an \underline{X} we can always replace that with basis function design matrix $\underline{\Phi}$.

Lasso (L1) regularization:

$$J_{\lambda}(\underline{w}) = \sum_{n=1}^N (y^{(n)} - f(\underline{x}^{(n)}; \underline{w}))^2 + \lambda \sum_{k=1}^K |w_k| \quad \dots \textcircled{2}$$

Still convex (unique minimum) but not "smooth" (differentiable). But other methods exist to solve (e.g. gradient descent instead of closed form method - later).

L_1 regularization has the effect of pushing weights to 0. This can be useful for interpreting data/a model (but be careful!)

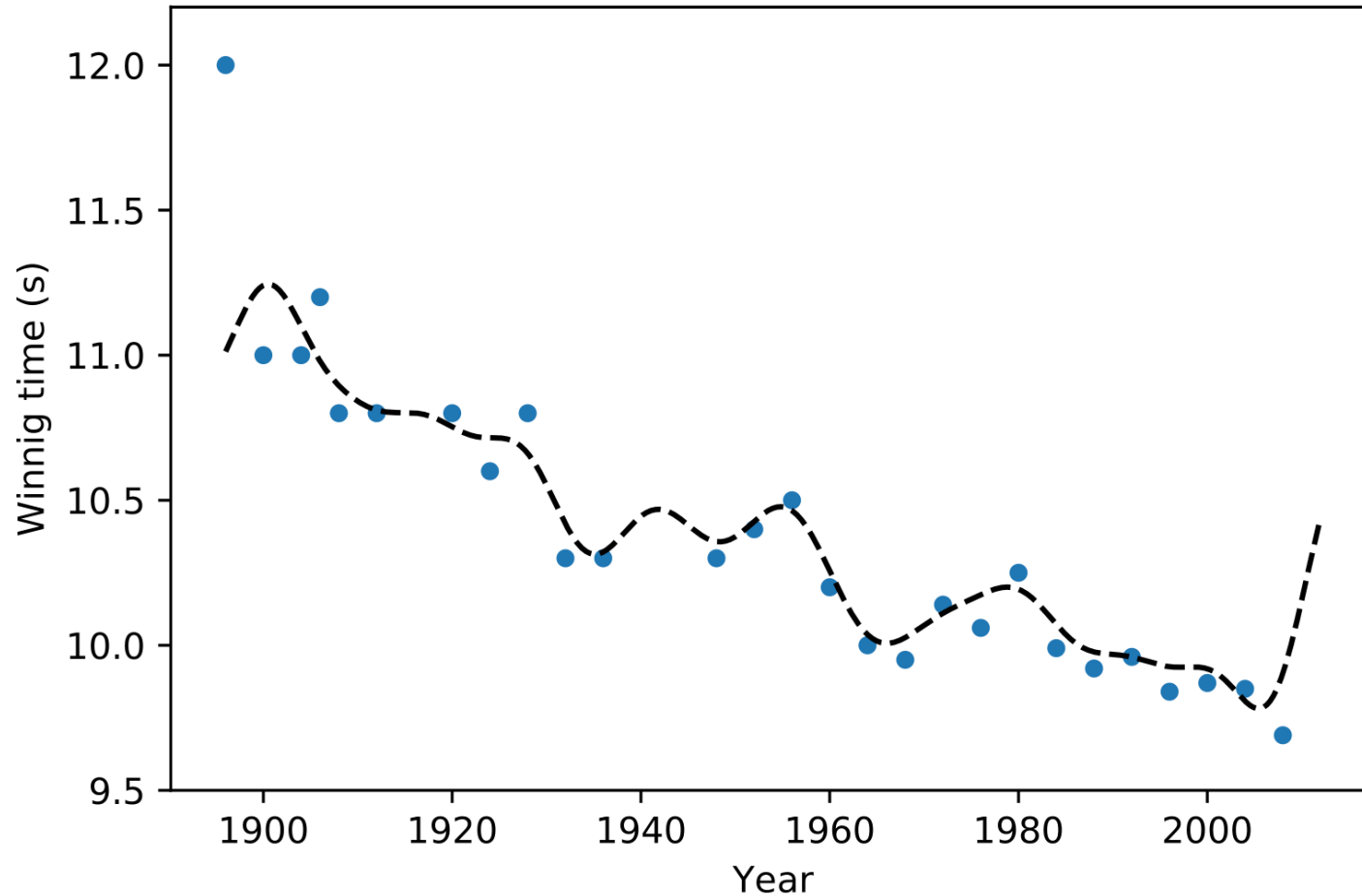
Why does L_1 do this but not L_2 ?

(Just intuitively from $\textcircled{1}$ and $\textcircled{2}$)



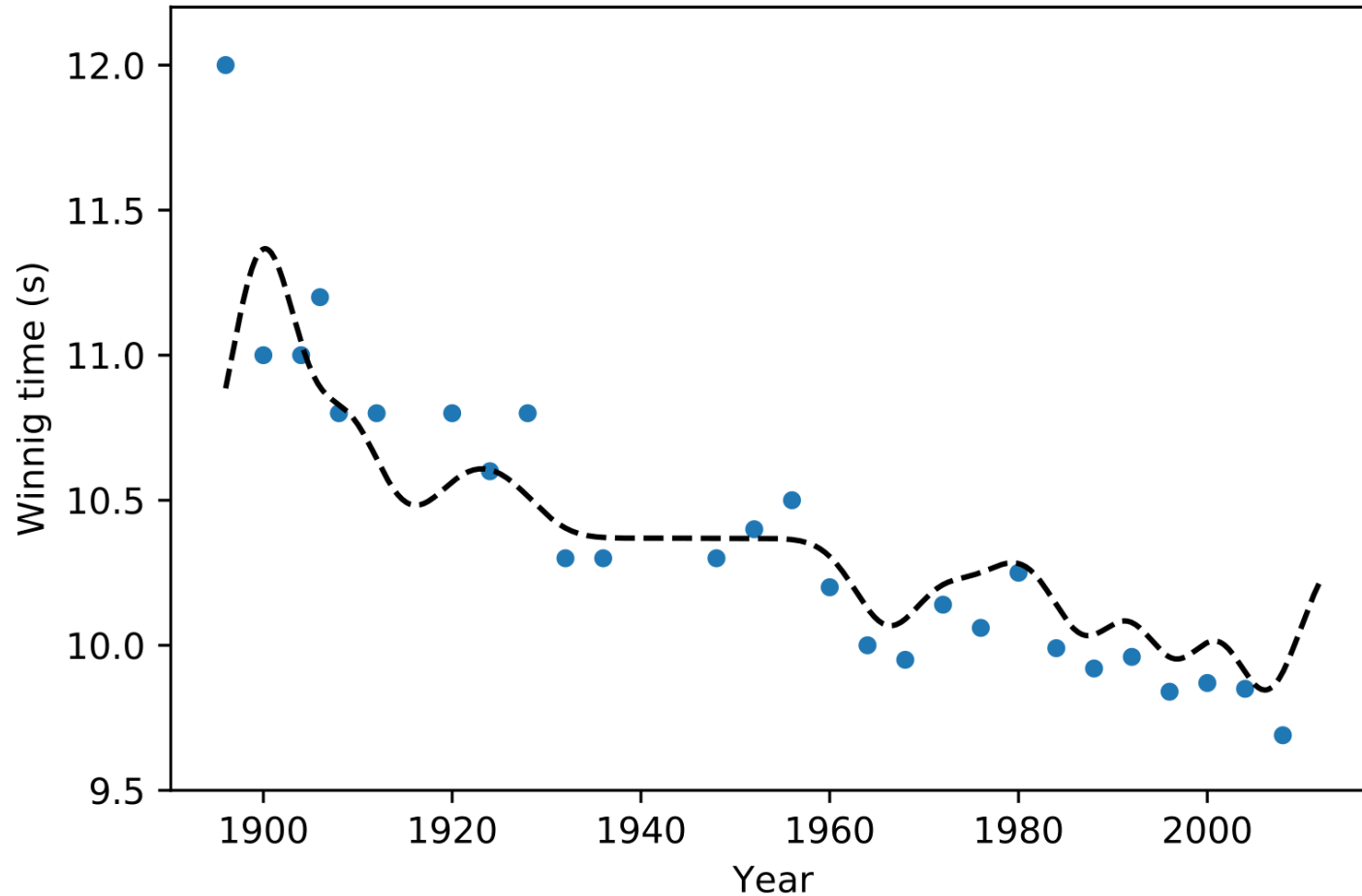
RBF with $c = [1900, 1901, \dots, 2000]$ and $h = 1$

Ridge regression with $\lambda = 1$: $\sum_{k=1}^K w_k^2 = 0.68$



RBF with $c = [1900, 1901, \dots, 2000]$ and $h = 1$

Lasso regression with $\lambda = 0.01$: $\sum_{k=1}^K w_k^2 = 1.52$



Lasso and ridge regression on diabetes data

