

Week 2, Lecture 3 - Fitting & Regression Redux, Regularization

Aaron Meyer

Outline

- ▶ Administrative Issues
- ▶ Fitting Regularization
 - ▶ Lasso
 - ▶ Ridge regression
 - ▶ Elastic net
- ▶ Some Examples

Based on slides from Rob Tibshirani.

The Bias-Variance Tradeoff

The Bias-Variance Tradeoff

Estimating β

- ▶ As usual, we assume the model:

$$y = f(\mathbf{z}) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ In regression analysis, our major goal is to come up with some good regression function

$$\hat{f}(\mathbf{z}) = \mathbf{z}^\top \hat{\beta}$$

- ▶ So far, we've been dealing with $\hat{\beta}^{ls}$, or the least squares solution:
 - ▶ $\hat{\beta}^{ls}$ has well known properties (e.g., Gauss-Markov, ML)
- ▶ But can we do better?

Choosing a good regression function

- ▶ Suppose we have an estimator

$$\hat{f}(\mathbf{z}) = \mathbf{z}^T \hat{\beta}$$

- ▶ To see if this is a good candidate, we can ask ourselves two questions:
 1. Is $\hat{\beta}$ close to the true β ?
 2. Will $\hat{f}(\mathbf{z})$ fit future observations well?
- ▶ These might have very different outcomes!!

Is $\hat{\beta}$ close to the true β ?

- ▶ To answer this question, we might consider the **mean squared error** of our estimate $\hat{\beta}$:
 - ▶ i.e., consider squared distance of $\hat{\beta}$ to the true β :

$$MSE(\hat{\beta}) = \mathbb{E} \left[\left\| \hat{\beta} - \beta \right\|^2 \right] = \mathbb{E} [(\hat{\beta} - \beta)^\top (\hat{\beta} - \beta)]$$

- ▶ **Example:** In least squares (LS), we know that:

$$\mathbb{E} [(\hat{\beta}^{ls} - \beta)^\top (\hat{\beta}^{ls} - \beta)] = \sigma^2 \text{tr} [(\mathbf{Z}^\top \mathbf{Z})^{-1}]$$

Will $\hat{f}(z)$ fit future observations well?

- ▶ Just because $\hat{f}(z)$ fits our data well, this doesn't mean that it will be a good fit to new data
- ▶ In fact, suppose that we take new measurements y'_i at the same \mathbf{z}_i 's:

$$(\mathbf{z}_1, \mathbf{y}'_1), (\mathbf{z}_2, \mathbf{y}'_2), \dots, (\mathbf{z}_n, \mathbf{y}'_n)$$

- ▶ So if $\hat{f}(\cdot)$ is a good model, then $\hat{f}(\mathbf{z}_i)$ should also be close to the new target y'_i
- ▶ This is the notion of **prediction error** (PE)

Prediction error and the bias-variance tradeoff

- ▶ So good estimators should, on average have, small prediction errors
- ▶ Let's consider the PE at a particular target point \mathbf{z}_0 :
 - ▶ $PE(\mathbf{z}_0) = \sigma_\epsilon^2 + Bias^2(\hat{f}(\mathbf{z}_0)) + Var(\hat{f}(\mathbf{z}_0))$
 - ▶ Not going to derive, but comes directly from previous definitions
- ▶ Such a decomposition is known as the **bias-variance tradeoff**
 - ▶ As model becomes more complex (more terms included), local structure/curvature is picked up
 - ▶ But coefficient estimates suffer from high variance as more terms are included in the model
- ▶ So introducing a little bias in our estimate for β might lead to a large decrease in variance, and hence a substantial decrease in PE

Depicting the bias-variance tradeoff

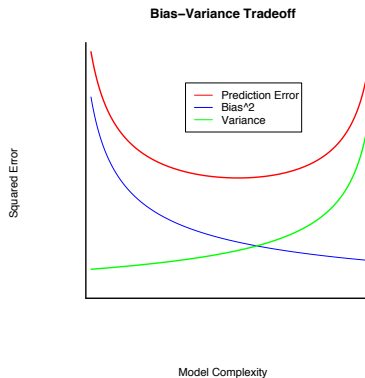


Figure: A graph depicting the bias-variance tradeoff.

Ridge Regression

Ridge Regression

Ridge regression as regularization

- ▶ If the β_j 's are unconstrained...
 - ▶ They can explode
 - ▶ And hence are susceptible to very high variance
- ▶ To control variance, we might **regularize** the coefficients
 - ▶ i.e., Might control how large the coefficients grow
- ▶ Might impose the ridge constraint (both equivalent):
 - ▶ minimize $\sum_{i=1}^n (y_i - \beta^\top \mathbf{z}_i)^2$ s.t. $\sum_{j=1}^p \beta_j^2 \leq t$
 - ▶ minimize $(\mathbf{y} - \mathbf{Z}\beta)^\top (\mathbf{y} - \mathbf{Z}\beta)$ s.t. $\sum_{j=1}^p \beta_j^2 \leq t$
- ▶ By convention (very important!):
 - ▶ \mathbf{Z} is assumed to be standardized (mean 0, unit variance)
 - ▶ \mathbf{y} is assumed to be centered

Ridge regression: l_2 -penalty

- ▶ Can write the ridge constraint as the following **penalized** residual sum of squares (PRSS):

$$\begin{aligned} PRSS(\boldsymbol{\beta})_{l_2} &= \sum_{i=1}^n (y_i - \mathbf{z}_i^\top \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_2^2 \end{aligned}$$

- ▶ Its solution may have smaller average PE than $\hat{\boldsymbol{\beta}}^{ls}$
- ▶ $PRSS(\boldsymbol{\beta})_{l_2}$ is convex, and hence has a unique solution
- ▶ Taking derivatives, we obtain:

$$\frac{\delta PRSS(\boldsymbol{\beta})_{l_2}}{\delta \boldsymbol{\beta}} = -2\mathbf{Z}^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) + 2\lambda \boldsymbol{\beta}$$

The ridge solutions

- ▶ The solution to $PRSS(\hat{\beta})_{l_2}$ is now seen to be:

$$\hat{\beta}_{\lambda}^{ridge} = (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^T \mathbf{y}$$

- ▶ Remember that \mathbf{Z} is standardized
 - ▶ \mathbf{y} is centered
- ▶ Solution is indexed by the tuning parameter λ (more on this later)
- ▶ Inclusion of λ makes problem non-singular even if $\mathbf{Z}^T \mathbf{Z}$ is not invertible
 - ▶ This was the original motivation for ridge regression (Hoerl and Kennard, 1970)

Tuning parameter λ

- ▶ Notice that the solution is indexed by the parameter λ
 - ▶ So for each λ , we have a solution
 - ▶ Hence, the λ 's trace out a path of solutions (see next page)
- ▶ λ is the shrinkage parameter
 - ▶ λ controls the size of the coefficients
 - ▶ λ controls amount of **regularization**
 - ▶ As λ decreases, we obtain the least squares solutions
 - ▶ As λ increases, we have $\hat{\beta}_{\lambda=0}^{ridge} = 0$ (intercept-only model)

Ridge coefficient paths

- ▶ The λ 's trace out a set of ridge solutions, as illustrated below

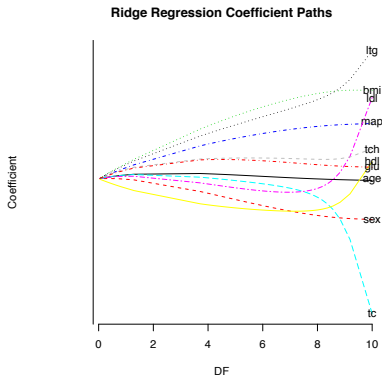


Figure: Ridge coefficient path for the diabetes data set found in the `lars` library in R.

Choosing λ

- ▶ Need disciplined way of selecting λ
- ▶ That is, we need to “tune” the value of λ
- ▶ In their original paper, Hoerl and Kennard introduced **ridge traces**:
 - ▶ Plot the components of $\hat{\beta}_\lambda^{ridge}$ against λ
 - ▶ Choose λ for which the coefficients are not rapidly changing and have “sensible” signs
 - ▶ No objective basis; heavily criticized by many
- ▶ Standard practice now is to use cross-validation (next lecture!)

A few notes on ridge regression

- ▶ The regularization decreases the degrees of freedom of the model
 - ▶ So you still cannot fit a model with more degrees of freedom than points
- ▶ This can be shown by examination of the smoother matrix
 - ▶ We won't do this—it's a complicated argument

How do we choose λ ?

- ▶ We need a disciplined way of choosing λ
- ▶ Obviously want to choose λ that minimizes the mean squared error
- ▶ Issue is part of the bigger problem of **model selection**

K-Fold Cross-Validation

- ▶ A common method to determine λ is K-fold cross-validation.
- ▶ **We will discuss this next lecture.**

Plot of CV errors and standard error bands

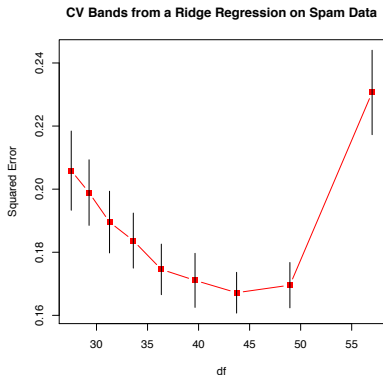


Figure: Cross validation errors from a ridge regression example on spam data.

The LASSO

The LASSO

The LASSO: l_1 penalty

- ▶ Tibshirani (*J of the Royal Stat Soc* 1996) introduced the **LASSO**: *least absolute shrinkage and selection operator*
- ▶ LASSO coefficients are the solutions to the l_1 optimization problem:

$$\text{minimize } (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) \text{ s.t. } \sum_{j=1}^p \|\beta_j\| \leq t$$

- ▶ This is equivalent to loss function:

$$\begin{aligned} PRSS(\boldsymbol{\beta})_{l_1} &= \sum_{i=1}^n (y_i - \mathbf{z}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p \|\beta_j\| \\ &= (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \end{aligned}$$

λ (or t) as a tuning parameter

- ▶ Again, we have a tuning parameter λ that controls the amount of regularization
- ▶ One-to-one correspondence with the threshold t :
 - ▶ recall the constraint:

$$\sum_{j=1}^p \|\beta_j\| \leq t$$

- ▶ Hence, have a “path” of solutions indexed by t
- ▶ If $t_0 = \sum_{j=1}^p \|\hat{\beta}_j^{ls}\|$ (equivalently, $\lambda = 0$), we obtain no shrinkage (and hence obtain the LS solutions as our solution)
- ▶ Often, the path of solutions is indexed by a fraction of shrinkage factor of t_0

Sparsity and exact zeros

- ▶ Often, we believe that many of the β_j 's should be 0
- ▶ Hence, we seek a set of **sparse solutions**
- ▶ Large enough λ (or small enough t) will set some coefficients exactly equal to 0!
 - ▶ So LASSO will perform model selection for us!

Computing the LASSO solution

- ▶ Unlike ridge regression, $\hat{\beta}_\lambda^{lasso}$ has no closed form λ
- ▶ Original implementation involves quadratic programming techniques from convex optimization
- ▶ But Efron *et al*, *Ann Statist*, 2004 proposed LARS (least angle regression), which computes the LASSO path efficiently
 - ▶ Interesting modification called is called forward stagewise
 - ▶ In many cases it is the same as the LASSO solution
 - ▶ Forward stagewise is easy to implement:
<https://www-stat.stanford.edu/~hastie/TALKS/nips2005.pdf>

Forward stagewise algorithm

- ▶ As usual, assume \mathbf{Z} is standardized and \mathbf{y} is centered
- ▶ Choose a small ϵ . The forward-stagewise algorithm then proceeds as follows:
 1. Start with initial residual $\mathbf{r} = \mathbf{y}$, and $\beta_1 = \beta_2 = \dots = \beta_p = 0$
 2. Find the predictor $\mathbf{Z}_j (j = 1, \dots, p)$ most correlated with \mathbf{r}
 3. Update $\beta_j = \beta_j + \delta_j$, where $\delta_j = \epsilon \cdot \text{sign}\langle \mathbf{r}, \mathbf{Z}_j \rangle = \epsilon \cdot \text{sign}(\mathbf{Z}_j^T \mathbf{r})$
 4. Set $\mathbf{r} = \mathbf{r} - \delta_j \mathbf{Z}_j$
 5. Repeat from step 2 many times

The LASSO, LARS, and Forward Stagewise paths

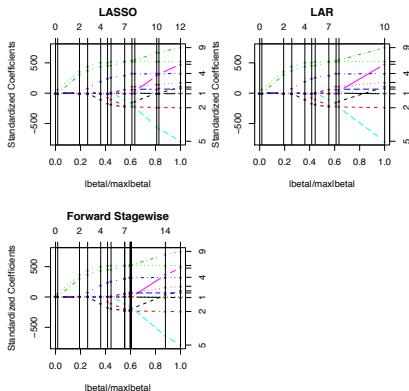


Figure: Comparison of the LASSO, LARS, and Forward Stagewise coefficient paths for the diabetes data set.

Comparing LS, Ridge, and the LASSO

- ▶ Even though $\mathbf{Z}^T\mathbf{Z}$ may not be of full rank, both ridge regression and the LASSO admit solutions
- ▶ We have a problem when $p \gg n$ (more predictor variables than observations)
 - ▶ But both ridge regression and the LASSO have solutions
 - ▶ Regularization tends to reduce prediction error

More comments on variable selection

- ▶ Now suppose $p \gg n$
- ▶ Of course, we would like a parsimonious model (Occam's Razor)
- ▶ Ridge regression produces coefficient values for each of the p -variables
- ▶ But because of its l_1 penalty, the LASSO will set many of the variables exactly equal to 0!
 - ▶ That is, the LASSO produces **sparse solutions**
- ▶ So LASSO takes care of model selection for us
 - ▶ And we can even see when variables jump into the model by looking at the LASSO path

Variants

- ▶ Zou and Hastie (2005) propose the **elastic net**, which is a convex combination of ridge and the LASSO
 - ▶ Paper asserts that the elastic net can improve error over LASSO
 - ▶ Still produces sparse solutions

High-dimensional data and underdetermined systems

- ▶ In many modern data analysis problems, we have $p \gg n$
 - ▶ These comprise “high-dimensional” problems
- ▶ When fitting the model $y = \mathbf{z}^T \boldsymbol{\beta}$, we can have many solutions
 - ▶ i.e., our system is *underdetermined*
- ▶ Reasonable to suppose that most of the coefficients are exactly equal to 0

But do these methods pick the right/true variables?

- ▶ Suppose that only S elements of β are non-zero
- ▶ Now suppose we had an “Oracle” that told us which components of the $\beta = (\beta_1, \beta_2, \dots, \beta_p)$ are truly non-zero
- ▶ Let β^* be the least squares estimate of this “ideal” estimator:
 - ▶ So β^* is 0 in every component that β is 0
 - ▶ The non-zero elements of β^* are computed by regressing y on only the S important covariates
- ▶ It turns out we get *really* close to this cheating solution without cheating!
 - ▶ Candes & Tao. *Ann Statist.* 2007.

The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity

Jordi Barretina^{1,2,3,*}, Giordano Caponigro^{4*}, Nicolas Stransky^{1*}, Kavitha Venkatesan^{4*}, Adam A. Margolin^{1†*}, Sungjoon Kim⁵, Christopher J. Wilson⁴, Joseph Lehar⁴, Gregory V. Kryukov¹, Dmitriy Sonkin⁴, Anupama Reddy⁴, Manway Liu⁴, Lauren Murray⁴, Michael F. Berger^{1*}, John E. Monahan⁴, Paula Morais⁴, Jodi Meltzer⁴, Adam Korejwa¹, Judit Jané-Valbuena^{1,2}, Felipa A. Mapa⁴, Joseph Thibault², Eva Bric-Furlong⁴, Pichai Raman⁴, Aaron Shipway³, Ingo H. Engels³, Jill Cheng⁶, Guoying K. Yu⁶, Jianjun Yu⁶, Peter Aspesi Jr⁴, Melanie de Silva⁴, Kalpana Jagtap⁴, Michael D. Jones⁴, Li Wang⁴, Charles Hattom³, Emanuele Palescandoli³, Supriya Gupta¹, Scott Mahan¹, Carrie Sougnez², Robert C. Onofrio¹, Ted Liefeld¹, Laura MacConaill², Wendy Winckler¹, Michael Reich¹, Nanxin Li², Jill P. Mesirov¹, Stacey B. Gabriel¹, Gad Getz¹, Kristin Ardlie³, Vivien Chan⁶, Vic E. Myer⁴, Barbara L. Weber¹, Jeff Porter⁴, Markus Warmuth⁴, Peter Finan⁴, Jennifer L. Harris², Matthew Meyerson^{1,2,3}, Todd R. Golub^{1,3,7,8}, Michael P. Morrissey^{4*}, William R. Sellers^{4*}, Robert Schlegel^{4*} & Levi A. Garraway^{1,2,3*}

The systematic translation of cancer genomic data into knowledge of tumour biology and therapeutic possibilities remains challenging. Such efforts should be greatly aided by robust preclinical model systems that reflect the genomic diversity of human cancers and for which detailed genetic and pharmacological annotation is available¹. Here we describe the Cancer Cell Line Encyclopedia (CCLE): a compilation of gene expression, chromosomal copy number and massively parallel sequencing data from 947 human cancer cell lines. When coupled with pharmacological profiles for 24 anticancer drugs across 479 of the cell lines, this collection allowed identification of genetic, lineage, and gene-expression-based predictors of drug sensitivity. In addition to known predictors, we found that plasma cell lineage correlated with sensitivity to IGF1 receptor inhibitors; *AHR* expression was associated with MEK inhibitor efficacy in *NRAS*-mutant lines; and *SLFN11* expression predicted sensitivity to topoisomerase inhibitors. Together, our results indicate that large, annotated cell-line collections may help to enable preclinical stratification schemata for anticancer agents. The generation of genetic predictions of drug response in the preclinical setting and their incorporation into cancer clinical trial design could speed the emergence of 'personalized' therapeutic regimens².

Human cancer cell lines represent a mainstay of tumour biology and drug discovery through facile experimental manipulation, global and

known cancer genes were assessed by mass spectrometric genotyping¹³ (Supplementary Table 2 and Supplementary Fig. 1). DNA copy number was measured using high-density single nucleotide polymorphism arrays (Affymetrix SNP 6.0, Supplementary Methods). Finally, messenger RNA expression levels were obtained for each of the lines using Affymetrix U133 plus 2.0 arrays. These data were also used to confirm cell line identities (Supplementary Methods and Supplementary Figs 2–4).

We next measured the genomic similarities by lineage between CCLE lines and primary tumours from Tumorscape¹⁴, expO, MILE and COSMIC data sets (Fig. 1b–d and Supplementary Methods). For most lineages, a strong positive correlation was observed in both chromosomal copy number and gene expression patterns (median correlation coefficients of 0.77, range = 0.52–0.94, $P < 10^{-15}$, for copy number, and 0.60, range = 0.29–0.77, $P < 10^{-15}$, for expression, respectively; Fig. 1b, c and Supplementary Tables 3 and 4), as has been described previously^{2,3,12}. A positive correlation was also observed for point mutation frequencies (median correlation coefficient = 0.71, range = -0.06–0.97, $P < 10^{-2}$ for all but 3 lineages; Supplementary Fig. 5), even when *TP53* was removed from the data set (median correlation coefficient = 0.64, range = -0.31–0.97, $P < 10^{-2}$ for all but 3 lineages; Fig. 1d and Supplementary Table 5). Thus, with relatively few exceptions (Supplementary Information), the CCLE may provide representative genetic proxies for primary tumours in many cancer types.

Example - Predicting Drug Response

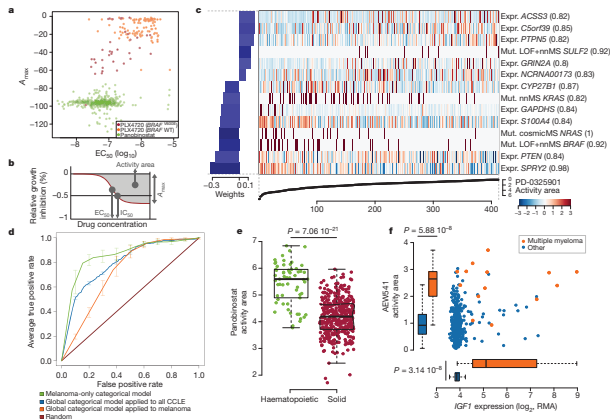


Figure 2 | Predictive modelling of pharmacological sensitivity using CCLE genomic data. **a, b**, Drug responses for panobinostat (green) and PLX4720 (orange/purple) represented by the high-concentration effect level (A_{max}) and transitional concentration (EC_{50}) for a sigmoidal fit to the response curve (**b**). **c**, Elastic net regression modelling of genomic features that predict sensitivity to PD-0325901. The bottom curve indicates drug response, measured as the area over the dose-response curve (activity area), for each cell line. The central heat map shows the CCLE features in the model (continuous z-score for expression and copy number, dark red for discrete mutation calls), across all cell lines (x axis). Bar plot (left): weight of the top predictive features for sensitivity (bottom) or insensitivity (top). Parentheses indicate features present in > 80% of models after bootstrapping. LOF, loss of function mutation; nmMS, non-neutral missense mutation (Supplementary Methods).

d, Specificity and sensitivity (receiver operating characteristic curves) of cross-validated categorical models predicting the response to a MEK inhibitor, PD-0325901 (activity area). Mean true positive rate and standard deviation ($n = 5$) are shown when models are built using all lines (global categorical model, in blue and orange), or within only melanoma lines (green). **e**, Activity area values for panobinostat between cell lines derived from haematopoietic ($n = 61$) and solid tumours ($n = 387$). The middle bar, median; box, inter-quartile range; bars extend to $1.5 \times$ the inter-quartile range. **f**, Distribution of activity area values for AEW541 relative to *IGF1* mRNA expression. Orange dots, multiple myeloma cell lines ($n = 14$); blue dots, cell lines from other tumour types ($n = 434$). Box-and-whisker plots show the activity area or mRNA expression distributions relative to each cell line type (line, median; box, inter-quartile range), with bars extending to $1.5 \times$ the inter-quartile range.

Example - Predicting Drug Response

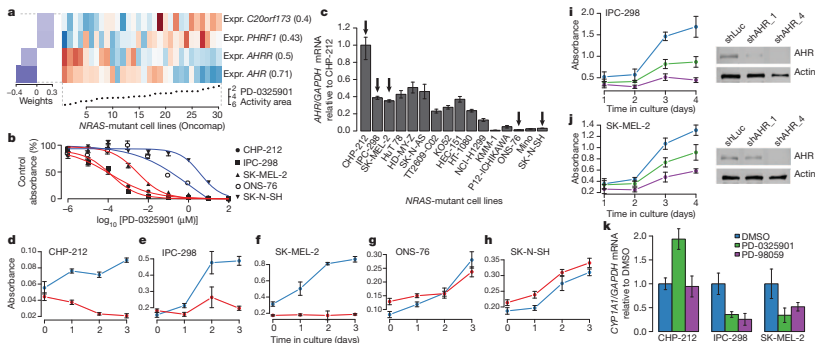


Figure 3 | *AHR* expression may denote a tumour dependency targeted by MEK inhibitors in *NRAS*-mutant cell lines. **a**, Predictive features for PD-0325901 sensitivity (using the 'varying baseline' activity area) in validated *NRAS*-mutant cell lines. **b**, Growth inhibition curves for *NRAS*-mutant cell lines expressing high (red) or low (blue) levels of *AHR* mRNA in the presence of the MEK inhibitor PD-0325901. **c**, Relative *AHR* mRNA expression across a panel of *NRAS*-mutant cell lines (arrows indicate cell lines where *AHR* dependency was analysed). **d-h**, Proliferation of *NRAS*-mutant cell lines displaying high (**d-f**) and low (**g, h**) *AHR* mRNA expression, after introduction of shRNAs against

AHR (red lines) or luciferase (blue lines). **i**, Left: proliferation of IPC-298 cells (high *AHR*) after introduction of additional shRNAs against *AHR* (shAHR_1 and shAHR_4; green and purple lines, respectively) or luciferase (control shLuc; blue line). Right: corresponding immunoblot analysis of *AHR* protein. **j**, Equivalent studies as in **i** using SK-MEL-2 cells (high *AHR*). **k**, Endogenous *CYP1A1* mRNA expression in the neuroblastoma line CHP-212 or the melanoma lines IPC-298 and SK-MEL-2 after exposure to vehicle (blue) or MEK inhibitors (PD-0325901, green or PD-98059, purple). Error bars indicate standard deviation between replicates, with $n = 12$ (**b**), $n = 3$ (**c**), $n = 6$ (**d-k**).

Example - Predicting Drug Response

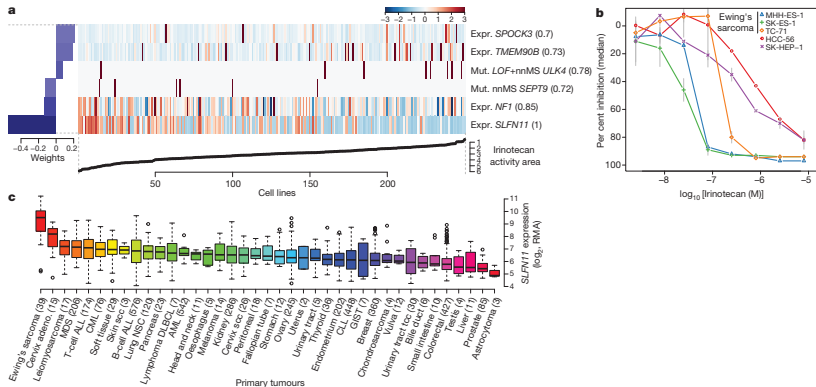


Figure 4 | Predicting sensitivity to topoisomerase I inhibitors. **a**, Elastic net regression analysis of genomic correlates of irinotecan sensitivity is shown for 250 cell lines. **b**, Dose-response curves for three Ewing's sarcoma cell lines (MSS-ES-1, SK-ES-1 and TC-71) and two control cell lines with low *SLFN11* expression (HCC-56 and SK-HEP-1). Grey vertical bars, standard deviation of

the mean growth inhibition ($n = 2$). **c**, *SLFN11* expression across 4,103 primary tumours. Box-and-whisker plots show the distribution of mRNA expression for each subtype, ordered by the median *SLFN11* expression level (line), the inter-quartile range (box) and up to $1.5\times$ the inter-quartile range (bars). Sample numbers (n) are indicated in parentheses.

Implementation

- ▶ The notebook can be found on the course website.

Further Reading

- ▶ Computer Age Statistical Inference, Chapter 16
- ▶ sklearn: Generalized Linear Models
- ▶ Candès E. and Tao T. The Dantzig selector: statistical estimation when p is much larger than n .