

Asymptotically Optimal Regularization in Smooth Parametric Models

Authored by:

Michael I. Jordan
Francis R. Bach
Percy S. Liang
Guillaume Bouchard

Abstract

Many types of regularization schemes have been employed in statistical learning, each one motivated by some assumption about the problem domain. In this paper, we present a unified asymptotic analysis of smooth regularizers, which allows us to see how the validity of these assumptions impacts the success of a particular regularizer. In addition, our analysis motivates an algorithm for optimizing regularization parameters, which in turn can be analyzed within our framework. We apply our analysis to several examples, including hybrid generative-discriminative learning and multi-task learning.

1 Paper Body

Many problems in machine learning and statistics involve the estimation of parameters from finite data. Although empirical risk minimization has favorable limiting properties, it is well known that this procedure can overfit on finite data. Hence, various forms of regularization have been employed to control this overfitting. Regularizers are usually chosen based on assumptions about the problem domain at hand. For example, in classification, we might use L2 regularization if we expect the data to be separable with a large margin. We might regularize with a generative model if we think it is roughly well-specified [7, 20, 15, 17]. In multi-task learning, we might penalize deviation between parameters across tasks if we believe the tasks to be similar [3, 12, 2, 13]. In each case, we would like (1) a procedure for choosing the parameters of the regularizer (for example, its strength) and (2) an analysis that shows the amount by which regularization reduces expected risk, expressed as a function of the compatibility between the regularizer and the problem domain. In this paper, we address these two points by developing an asymptotic analysis of smooth regularizers for parametric problems. The key idea is to derive a second-order

Taylor approximation of the expected risk, yielding a simple and interpretable quadratic form which can be directly minimized with respect to the regularization parameters. We first develop the general theory (Section 2) and then apply it to some examples of common regularizers used in practice (Section 3).

2

General theory

We use uppercase letters (e.g., L , R , Z) to denote random variables and script letters (e.g., \mathcal{L} , \mathcal{R} , \mathcal{I}) to denote constant limits ... of random variables. For a θ -parametrized differentiable function $f(\theta; \mathbf{z})$, let f' , f'' , and f''' denote the first, second and third derivatives of f with respect to θ , and let $\partial f(\theta; \mathbf{z})$ denote the derivative with respect to \mathbf{z} . Let $\mathbf{X}_n = \text{Op}(\mathbf{z}^n)$ denote a sequence of 1

P

random variables for which $\mathbf{z}^n \mathbf{X}_n$ is bounded in probability. Let $\mathbf{X}_n \rightarrow \mathbf{X}$ denote convergence in probability. For a vector \mathbf{v} , let $\mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|^2$. Expectation and variance operators are denoted as $E[\cdot]$ and $V[\cdot]$, respectively. 2.1

Setup

We are given a loss function $\ell(\theta; \mathbf{z})$ parametrized by $\theta \in \mathbb{R}^d$ (e.g., $\ell((\mathbf{x}, \mathbf{y}); \theta) = \frac{1}{2} (\mathbf{y} - \mathbf{x}^T \theta)^2$ for linear regression). Our goal is to minimize the expected risk, def

$$\theta^* = \arg\min_{\theta} L(\theta),$$

def

$$L(\theta) = E_{\mathbf{Z} \sim p} [\ell(\mathbf{Z}; \theta)],$$

(1)

\mathbb{R}^d

which averages the loss over some true data generating distribution $p(\mathbf{Z})$. We do not have access to p , but instead receive a sample of n i.i.d. data points $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ drawn from p . The standard unregularized estimator minimizes the empirical risk: $\hat{\theta}_n = \arg\min_{\theta} L_n(\theta)$, $L_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{Z}_i; \theta)$. (2) Although $\hat{\theta}_n$ is consistent (that is, it converges in probability to θ^*) under relatively weak conditions, it is well known that regularization can improve performance substantially for finite n . Let $R_n(\theta, \lambda)$ be a (possibly data-dependent) regularization function, where $\lambda \in \mathbb{R}_+$ are the regularization parameters. For linear regression, we might use squared regularization ($R_n(\theta, \lambda) = \frac{\lambda}{2} \|\theta\|^2$), where $\lambda \in \mathbb{R}_+$ determines the strength. Define the regularized estimator as follows: def $\hat{\theta}_n(\lambda) = \arg\min_{\theta} L_n(\theta) + R_n(\theta, \lambda)$. (3) \mathbb{R}_+

\mathbb{R}_+

The goal of this paper is to choose good values of λ and analyze the subsequent impact on performance. Specifically, we wish to minimize the relative risk: def $L_n(\lambda) = E_{\mathbf{Z}_1, \dots, \mathbf{Z}_n \sim p} [L(\hat{\theta}_n(\lambda)) - L(\theta^*)]$, (4) 1

$\frac{1}{n}$

$\frac{1}{n}$

$\frac{1}{n}$

which is the difference in risk (averaged over the training data) between the regularized and unregularized estimators; $L_n(\lambda) \rightarrow 0$ is desirable. Clearly, $\arg\min_{\lambda} L_n(\lambda)$ is the optimal regularization parameter. However, it is difficult

λ Mahalanobis metric given by L . Note that the squared regularizer bias is always positive: it always increases the risk by an amount which depends on how "wrong" the regularizer is.

λ L^{-1} }: The key quantity is $R(\lambda)$, λ Variance reduction provided by the regularizer $\text{tr}\{I - L^{-1} R(\lambda)\}$ the Hessian of the regularizer, whose impact on the relative risk is channeled through L^{-1} and λ I . For convex regularizers, $R(\lambda) \succeq 0$, so we always improve the stability of the estimate by regularizing. Furthermore, if the loss is the negative log-likelihood and our model is well-specified (that is, $p^*(z) = \exp\{\langle z; \theta^* \rangle\}$), then $I = L$ by the first Bartlett identity [4], and the variance λ L^{-1} }. reduction term simplifies to $\text{tr}\{R(\lambda)\}$. λ $\text{tr}\{I - R(\lambda)L^{-1}\}$: Alignment between regularizer bias and unregularized estimator bias $2B \preceq R(\lambda)$ The alignment has two parts, the first of which is nonzero only for non-linear models and the second of which is nonzero only when the regularizer depends on the training data. The unregularized λ estimator errs in direction B ; we can reduce the risk if the regularizer bias $R(\lambda)$ helps correct for the λ λ estimator bias ($B \preceq R(\lambda) \preceq 0$). The second part carries the same intuition: the risk is reduced when the random regularizer compensates for the loss ($\text{tr}\{I - R(\lambda)L^{-1}\} \leq 0$). 2.4

Oracle regularizer

The principal advantage of having a simple expression for $L(\lambda)$ is that we can minimize it with λ def respect to λ . Let $\lambda^* = \text{argmin}_{\lambda} L(\lambda)$ and call λ^* the oracle estimator. We have a closed form for λ^* in the important special case that the regularization parameter λ is the strength of the regularizer: Corollary 1 (Oracle regularization strength). If $R_n(\lambda, \lambda) = n^{-1} r(\lambda)$ for some $r(\lambda)$, then $\lambda^* = \text{argmin}_{\lambda} L(\lambda) = \lambda^*$

$$\begin{aligned} \text{tr}\{I - L^{-1} r L^{-1}\} + 2B \preceq r \text{ def } C1 = , C2 r \preceq L^{-1} r \\ L(\lambda^*) = \lambda^* \\ C12 \preceq 2C2 \\ (7) \end{aligned}$$

Proof. (6) is a quadratic in λ ; solve by differentiation. Compute $L(\lambda^*)$ by substitution. In general, λ^* will depend on λ and hence is not computable from data; Section 2.5 will remedy this. Nevertheless, the oracle regularizer provides an upper bound on performance and some insight into the relevant quantities that make a regularizer useful. Note $L(\lambda^*) \succeq 0$, since optimizing λ must be no worse than not regularizing since $L(0) = 0$. But what might be surprising at first is that the oracle regularization parameter λ^* can be negative 3

Estimator	Notation	Relative risk
UNREGULARIZED	$\lambda^* \succeq 0$	
ORACLE	$\lambda^* \preceq n^{-1} L(\lambda^*)$	
PLUGIN	$\lambda^* \preceq n^{-1} L(\lambda^*)$ (1)	
ORACLE PLUGIN	$\lambda^* \preceq n^{-1} L(\lambda^*)$ (???)	

Table 1: Notation for the various estimators and their relative risks. (corresponding to "anti-regularization"). But if helps ($\lambda^* \preceq 0$ and $L(\lambda) \preceq 0$ for $0 \preceq \lambda \preceq 2\lambda^*$). 2.5

$\lambda^* L(\lambda^*) \preceq$
 $= \lambda^* C1 \preceq 0$, then (positive) regularization

Plugin regularizer

While the oracle regularizer $R_n(\eta, \cdot)$ given by (7) is asymptotically optimal, η depends on the θ unknown θ , so $\hat{\eta}_n$ is actually not implementable. In this section, we develop the plugin regularizer \hat{R}_n as a way to avoid this dependence. The key idea is to substitute η with an estimate $\hat{\eta} = \hat{\eta}_n + \frac{1}{n} \text{def } \hat{\eta}_n(\cdot, \cdot)$, where $\hat{\eta}_n = \text{Op}(n^{-1/2})$. We then use the plugin estimator $\hat{\eta}_n = \arg\min_{\eta} L_n(\eta) + R_n(\hat{\eta}_n)$

$\hat{\eta}_n$ How well does this plugin estimator work, that is, what is its relative risk $E[L(\hat{\eta}_n) - L(\eta_0)]$? $\hat{\eta}_n$ and apply Theorem 1 because $L(\eta)$ can only be applied to nonrandom arguments. However, we can still leverage existing machinery by defining a new plugin $\hat{\eta}_n(\cdot, \cdot)$ with regularization parameter $\hat{\eta}_n \in \mathbb{R}$. Henceforth, the regularizer $R_n(\eta, \cdot) = \hat{R}_n(\eta, \cdot)$ superscript $\hat{\cdot}$ will denote quantities concerning the plugin regularizer. The corresponding estimator $\hat{\eta}_n = \arg\min_{\eta} L_n(\eta) + \hat{R}_n(\hat{\eta}_n, \cdot)$ has relative risk $L^*(\hat{\eta}_n) = E[L(\hat{\eta}_n) - L(\eta_0)]$. The key $\hat{\eta}_n$ identity is $\hat{\eta}_n = \hat{\eta}_n + \frac{1}{n}$, which means the asymptotic risk of the plugin estimator $\hat{\eta}_n$ is simply $L^*(1)$.

We could try to squeeze more out of the plugin regularizer by further optimizing $\hat{\eta}_n$ according to $\hat{\eta}_n = \arg\min_{\eta} L(\hat{\eta}_n)$ and use the oracle plugin estimator $\hat{\eta}_n$ rather than just using $\hat{\eta}_n = 1$. In general, this is not useful since $\hat{\eta}_n$ might depend on θ , and the whole point of plugin is to remove this dependence. However, in a fortuitous turn of events, for some linear models $\hat{\eta}_n$ (Sections 3.1 and 3.4), $\hat{\eta}_n$ is in fact independent of θ , and so $\hat{\eta}_n$ is actually implementable. Table 1 summarizes all the estimators we have discussed. The following theorem relates the risks of all estimators we have considered (see [16] for the proof): Theorem 2 (Relative risk of plugin). The relative risk of the plugin estimator is $L^*(1) = L(\eta_0) + E$, where $E = \lim_{n \rightarrow \infty} n E[\text{tr}\{L'(\eta_0) \hat{\eta}_n - L'(\eta_0) \eta_0\}]$. If $R_n(\eta)$ is linear in η , then the relative risk of the oracle plugin estimator is $L^*(\eta_0) = L^*(1) +$

$$\begin{aligned} & E[2L(\eta_0)] \\ & \text{with } \eta_0 = 1 + \\ & E[2L(\eta_0)] . \end{aligned}$$

Note that the sign of E depends on the nature of the error η_n , so P LUGIN could be either better or worse than O RACLE. On the other hand, O RACLE P LUGIN is always better than P LUGIN. We can get a simpler expression for E if we know more about η_n (see [16] for the proof): $\eta_n = f(\eta_0)$, then the Theorem 3. Suppose $\eta_n = f(\eta_0)$ for some differentiable $f: \mathbb{R}^d \rightarrow \mathbb{R}^b$. If η_n η_0 results of Theorem 2 hold with $E = \text{tr}\{L'(\eta_0) f'(\eta_0) L'(\eta_0)\}$.

3

Examples

In this section, we apply our results from Section 2 to specific problems. Having made all the asymptotic derivations in the general setting, we now only need to make a few straightforward calculations to obtain the asymptotic relative risks and regularization parameters for a given problem. We first explore two classical examples from statistics (Sections 3.1 and 3.2) to get some intuition

for the theory. Then we consider two important examples in machine learning (Sections 3.3 and 3.4).

3.1 Estimation of normal means

Assume that data are generated from a multivariate normal distribution with d independent components ($p = N(\mu, I)$). We use the negative log-likelihood as the loss function: $\ell(x; \mu) = \frac{1}{2} \|x - \mu\|^2$, so the model is well-specified.

In his seminal 1961 paper [14], Stein showed that, surprisingly, the standard empirical risk minimizer $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is beaten by the James-Stein estimator $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n+2} \frac{\sum_{i=1}^n \langle X_i, \hat{\mu}_n \rangle X_i}{\|\hat{\mu}_n\|^2}$ in the sense JS that $E[\|\hat{\mu}_n\|^2] \leq E[\|\hat{\mu}_n^{\text{JS}}\|^2]$ for all n and $d \geq 3$. We will show that the James-Stein estimator is essentially equivalent to O RACLE P LUGIN with quadratic regularization ($r(\mu) = \frac{1}{2} \|\mu\|^2$). Let $L = I$, $B = 0$, $r = \frac{1}{2} \|\cdot\|^2$, and $r' = I$. By (7), the oracle regularization First compute $L\mu = \frac{1}{n} \sum_{i=1}^n X_i$, $2d$ weight is $\frac{1}{n+2} \|\mu\|^2$, which yields a relative risk of $L(\mu) = \frac{1}{n+2} \|\mu\|^2$. Now let us derive P LUGIN (Section 2.5). We have $f(\mu) = \frac{1}{2} \|\mu\|^2$

$$\begin{aligned} & \frac{1}{n+2} \|\mu\|^2 \\ & \text{and } f'(\mu) = \\ & \frac{1}{n+2} \|\mu\|^2. \end{aligned}$$

By Theorems 2
 $d(d+2) \leq 2k \leq 2n$

and $L(1) = 0$. Note that since $E \geq 0$, P LUGIN is always (asymptotically) worse than O RACLE but better than U NREGULARIZED if $d \geq 4$. To get O RACLE P LUGIN, compute $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (note that this doesn't depend on μ), which results in $\frac{1}{n} \sum_{i=1}^n X_i$

$\frac{1}{n+2} \|\mu\|^2 \leq \frac{1}{n} \sum_{i=1}^n \langle X_i, \hat{\mu}_n \rangle \leq \frac{1}{n} \sum_{i=1}^n \|X_i\|^2$. By Theorem 2, its relative risk is $L(\mu) = \frac{1}{n+2} \|\mu\|^2$, which offers a small improvement over P LUGIN (and is superior to U NREGULARIZED when $d \geq 2$).

Note that the O RACLE P LUGIN and P LUGIN are adaptive: We regularize more or less depending on whether our preliminary estimate of X

$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$, which differs from μ by a very small amount: $\|\hat{\mu}_n - \mu\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|X_i - \mu\|^2$. O RACLE P LUGIN has the added benefit that it always shrinks towards zero by an amount between 0 and 1, whereas JAMES STEIN can overshoot. Empirically, we found that O RACLE P LUGIN generally had a lower expected risk than JAMES STEIN when $k \leq 1$, but JAMES STEIN was better when $k > 1$.

Binomial estimation

Consider the estimation of θ , the log-odds of a coin coming up heads. We use the negative loglikelihood loss $\ell(x; \theta) = -x\theta + \log(1 + e^\theta)$, where $x \in \{0, 1\}$ is the outcome of the coin. This example serves to provide intuition for the bias B appearing in (6), which is typically ignored in first-order asymptotics or is zero (for linear models). Consider a regularizer $r(\theta) = \frac{1}{2} (\theta + 2 \log(1 + e^\theta))$, which corresponds to a $\text{Beta}(\frac{1}{2}, \frac{1}{2})$ prior. Choosing θ has been studied

extensively in statistics. Some common choices are the Haldane prior ($\alpha = 0$), the reference (Jeffreys) prior ($\alpha = 1$), the uniform prior ($\alpha = 2$), and Laplace smoothing ($\alpha = 4$). We will choose α to minimize expected risk adaptively based on data. ... Define $\alpha = 1 + e^{1/\alpha}$, $v = \alpha(1 - \alpha)$, and $b = \alpha - 1/2$. Then compute $L = v$, $L = \alpha vb$, $r = b$, $r = v$, $B = \alpha v - b$. O RACLE corresponds to $\alpha = 2 + bv^2$. Note that $\alpha \geq 0$, so again (positive) regularization always helps.

We can compute the difference between O RACLE and P LUGIN: $E = 2\alpha v b^2$. If $b \geq 4$, $E \geq 0$, which means that P LUGIN is worse; otherwise P LUGIN is actually better. Even when P LUGIN is worse than O RACLE, P LUGIN is still better than U NREGULARIZED, which can be verified by checking that $L(1) = 25 vb^2 - 2v - 1 b^2 \geq 0$ for all α .

3.3

Hybrid generative-discriminative learning

In prediction tasks, we wish to learn a mapping from some input $x \in X$ to an output $y \in Y$. A common approach is to use probabilistic models defined by exponential families, which is defined by a vector of sufficient statistics (features) $\phi(x, y) \in \mathbb{R}^d$ and an accompanying vector of parameters $\theta \in \mathbb{R}^d$. These features can be used to define a generative model (8) or a discriminative model (9): $Z p_\theta(x, y) = \exp\{\phi(x, y) \cdot \theta\}$, $A(\theta) = \log \int \exp\{\phi(x, y) \cdot \theta\} dy dx$, (8) $X Y Z p_\theta(y - x) = \exp\{\phi(x, y) \cdot \theta\}$, $A(\theta; x) = \log \int \exp\{\phi(x, y) \cdot \theta\} dy$. (9) Y

5

Misspecification 0% 5% 50%

$\text{tr}\{I^{-1} \text{Cov}(\phi) \} = 5.38 \ 13.8$

$2B(\theta) = 0 \ -0.073 \ -1.0$

$\text{tr}\{(\text{Cov}(\phi) \text{Cov}(\phi))^2 \} = 0 \ 0.00098 \ 0.034$

$\text{tr}\{(\text{Cov}(\phi) \text{Cov}(\phi))^3 \} = 310 \ 230$

$L(\theta) = -0.65 \ -48 \ -808$

Table 2: The oracle regularizer for the hybrid generative-discriminative estimator. As misspecification increases, we regularize less, but the relative risk is reduced more (due to more variance reduction). Given these definitions, we can either use a generative estimator $\hat{\theta}_{\text{ngen}} = \arg\min_{\theta} G_n(\theta)$, where $P_n G_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(x, y)$ or a discriminative estimator $\hat{\theta}_{\text{ndis}} = \arg\min_{\theta} D_n(\theta)$, where $P_n D_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(y - x)$.

There has been a flurry of work on combining generative and discriminative learning [7, 20, 15, 18, 17]. [17] showed that if the generative model is well-specified ($p_\theta(x, y) = p^*(x, y)$), then the generative estimator is better in the sense that $L(\hat{\theta}_{\text{ngen}}) \leq L(\hat{\theta}_{\text{ndis}}) \leq nc + O_p(n^{-1/2})$ for some $c \geq 0$; if the model is misspecified, the discriminative estimator is asymptotically better. To create a hybrid estimator, let us treat the discriminative and generative objectives as the empirical risk and the regularizer, respectively, so $\ell((x, y); \theta) = -\log p_\theta(y - x)$, so $L_n = D_n$ and $R_n(\theta, \theta) = n G_n(\theta)$. As $n \rightarrow \infty$, the discriminative objective dominates as desired. Our approach generalizes the analysis of [6], which applies only to unbiased estimators for conditionally well-specified models. By moment-generating properties of the exponential family,

we arrive at the following quantidef def λ ties (write λ for $\lambda(X, Y)$): $L = \text{vx} = \text{Ep}_{\lambda} (X) [\text{Vp}_{\lambda} (Y - X) [\lambda - X]]$, $R(\lambda) = \lambda(\lambda, \lambda, \lambda_{xy}) = \text{def } \lambda(\text{Ep}_{\lambda} (X, Y) [\lambda], \lambda, \lambda_{xy})$, and $R(\lambda) = \lambda_v = \lambda \text{Vp}_{\lambda} (X, Y) [\lambda]$. The oracle regularization parameter is then $\lambda^* =$

$$\frac{\text{tr}\{I \text{“ vx?1 vx?1 } \} + 2B \lambda(\lambda, \lambda, \lambda_{xy})}{\text{tr}\{I \text{“ vx?1 } \} \cdot \text{tr}\{(\lambda, \lambda, \lambda_{xy}) \text{ vx?1 } \}} \quad (10)$$

The sign and magnitude of λ^* provides insight into how generative regularization improves prediction as a function of the model and problem: Specifically, a large positive λ^* suggests regularization is helpful. To simplify, assume that the discriminative model is well-specified, that is, $p^*(y - x) = p^{**}(y - x)$ (note that the generative model could still be misspecified). In this case, $\lambda^* \text{I}^* = \text{vx}$, and so the numerator reduces to $\text{tr}\{(v - \text{vx}) \text{vx?1 } \} + 2B \lambda(\lambda, \lambda, \lambda_{xy})$. $\text{I}^* = L$. Since $v - \text{vx}$ (the key fact used in [17]), the variance reduction (plus the random alignment term from I^*) is always non-negative with magnitude equal to the fraction of missing information provided by the generative model. There is still the non-random alignment term $2B \lambda(\lambda, \lambda, \lambda_{xy})$, whose sign depends on the problem. Finally, the denominator (always positive) affects the optimal magnitude of the regularization. If the generative model is almost well-specified, λ^* will be close to λ_{xy} , and the regularizer should be trusted more (large λ^*). Since our analysis is local, misspecification (how much $p^{**}(x, y)$ deviates from $p^*(x, y)$) is measured by a Mahalanobis distance between λ and λ_{xy} , rather than something more stringent and global like KL-divergence. An empirical example To provide some concrete intuition, we investigated the oracle regularizer for a synthetic binary classification problem of predicting $y \in \{0, 1\}$ from $x \in \{0, 1\}^k$. Using features $\lambda(x, y) = (I[y = 0]x_i, I[y = 1]x_i)_i$ defines the corresponding generative (Naive Bayes) λ and discriminative (logistic regression) estimators. We set $k = 5$, $\lambda = (10, \lambda, \lambda, \lambda, \lambda)$, $\lambda = (10, \lambda, \lambda, \lambda, \lambda)$, λ and $p(x, y) = (1 - \lambda)p^{**}(x, y) + \lambda p^{**}(y)p^{**}(x_1 - y)I[x_1 = \lambda] = \lambda k$. The amount of misspecification is controlled by $0 \leq \lambda \leq 1$, the fraction of examples whose features are perfectly correlated. Table 2 shows how the oracle regularizer changes with λ . As λ increases, λ^* decreases (we regularize less) as expected. But perhaps surprisingly, the relative risk is reduced with more misspecification; this is due to the fact that the variance reduction term increases and has a quadratic effect on $L(\lambda)$. Figure 1(a) shows the relative risk $\text{Ln}(\lambda)$ for various values of λ . The vertical line corresponds to λ^* , which was computed numerically by sampling. Note that the minimum of the curves

$(\arg\min_{\lambda} \text{Ln}(\lambda))$, the desired quantity, is quite close to λ^* and approaches λ^* as n increases, which empirically justifies our asymptotic approximations. Unlabeled data One of the main advantages of having a generative model is that we can leverage unlabeled examples by marginalizing out their hidden outputs. Specifically, suppose we have m i.i.d. unlabeled examples $X_{n+1}, \dots, X_{n+m} \sim p(x)$, with $m \rightarrow \infty$ as $n \rightarrow \infty$. Define the m unlabeled regularizer as $R_n(\lambda, \lambda) = \frac{1}{m} \sum_{i=1}^m \log p^*(X_{n+i})$. We can compute $R^* = \lambda^*$ using the stationary conditions of the loss function at λ^* . Also, $\lambda^* = v - \text{vx}$,

and $\lambda_r = 0$ (the regularizer doesn't depend on the labeled data). If the model is R conditionally well-specified, we can verify that the oracle regularization parameter λ^* is the same as if we had regularized with G_n . This equivalence suggests that the dominant concern asymptotically is developing an adequate generative model with small bias and not exactly how it is used in learning.

3.4 Multi-task regression

The intuition behind multi-task learning is to share statistical strength between tasks [3, 12, 2, 13]. Suppose we have K regression tasks. For each task $k = 1, \dots, K$, we generate each data point (x_k, y_k) . We can treat this $i = 1, \dots, n$ independently as follows: $X_{ik} \sim p_i(X_{ik})$ and $Y_{ik} \sim N(X_{ik}^T \beta_k, \sigma_k^2)$ as a single task problem by concatenating the vectors for all the tasks: $X_i = (X_{i1}, \dots, X_{iK})^T \in \mathbb{R}^{dK}$, $Y = (Y_1, \dots, Y_K)^T \in \mathbb{R}^K$, and $\beta = (\beta_1, \dots, \beta_K)^T \in \mathbb{R}^{dK}$. It will also be useful to represent β by the matrix $\beta = (\beta_1, \dots, \beta_K) \in \mathbb{R}^{d \times K}$. The loss function is $\ell((x, y), \beta) = \frac{1}{2} \sum_{k=1}^K (y_k - x_k^T \beta_k)^2$. Assume the model is conditionally well-specified. We would like to be flexible in case some tasks are more related than others, so let us define a positive definite matrix $\Lambda \in \mathbb{R}^{K \times K}$ of inter-task affinities and use the quadratic regularizer: $r(\beta) = \frac{1}{2} \beta^T \Lambda \beta$, which implies that $I = L = \Lambda K$. For simplicity, assume $E X_i X_i^T = I$. Most of the computations that follow parallel those of Section 3.1, only extended to matrices. Substituting the relevant quantities into (6) yields the relative risk: $L(\beta) = \frac{1}{2} \text{tr}\{\beta^T \Lambda \beta\} + \frac{1}{2} \text{tr}\{\beta^T \Lambda \beta\}$. Optimizing with respect to β produces the oracle regularization parameter $\beta^* = d(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}$; its associated relative risk $L(\beta^*) = \frac{1}{2} d \text{tr}\{\Lambda^{-1}\}$. To analyze P LUGIN, first compute $f = d(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}$. However, the relative risk of P LUGIN is increases the asymptotic risk by $E = 2 \text{tr}\{(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}\}$ still favorable when $d \geq 4$, as $L(1) = \frac{1}{2} d(d-4) \text{tr}\{(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}\} \geq 0$ for $d \geq 4$.

We can do slightly better using O RACLE P LUGIN ($\lambda = 1/d^2$), which results in a relative risk of $L(\lambda) = \frac{1}{2} d(d-2) \text{tr}\{(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}\}$. For comparison, if we had solved the K regression tasks completely independently with K independent regularization parameters, our relative risk would have been $\frac{1}{2} d(d-2) \sum_{k=1}^K \text{tr}\{(\Lambda^{-1} \Lambda + \Lambda^{-1})^{-1} \Lambda^{-1}\}$ (following similar but simpler computations). We now compare joint versus independent regularization. Let $A = \Lambda^{-1} \Lambda$ with eigendecomposition $A = U D U^T$. The difference in relative risks between joint and independent regularization is $\Delta = \frac{1}{2} d(d-2) \sum_{k=1}^K (\lambda_k - 1)$ ($\lambda_k \geq 0$ means joint regularization is better). The gap Δ between joint and independent regularization is large when the tasks are non-trivial but similar (λ_k are close to 1, but $\lambda_k - 1$ is large). In that case, D_{kk} is quite large for $k \geq 1$, but all the A_{kk} are small. MHC-I binding prediction We evaluated our multitask regularization method on the IEDB MHC-I peptide binding dataset created by [19] and used by [13]. The goal here is to predict the binding affinity (represented by $\log IC_{50}$) of a MHC-I molecule given its amino-acid sequence (represented by a vector of binary features, reduced to a 20-dimensional real vector using SVD). We created five regression tasks corresponding to the five most common MHC-I molecules. We compared four estimators: UNREGULARIZED, DIAGONAL CV ($\lambda = cI$), UNIFORM CV (using the same task-affinity

for all pairs of tasks with $\gamma = c(1/\sqrt{n} + 10^{-5})$, and P LUGIN CV ($\gamma = \gamma_i/\sqrt{n}$), where c was chosen by cross-validation.³ Figure 1 shows the results averaged over $\text{cd}(\sqrt{n}/n) \cdot 3$

We performed three-fold cross-validation to select c from 21 candidates in $[10^{-5}, 10^5]$.

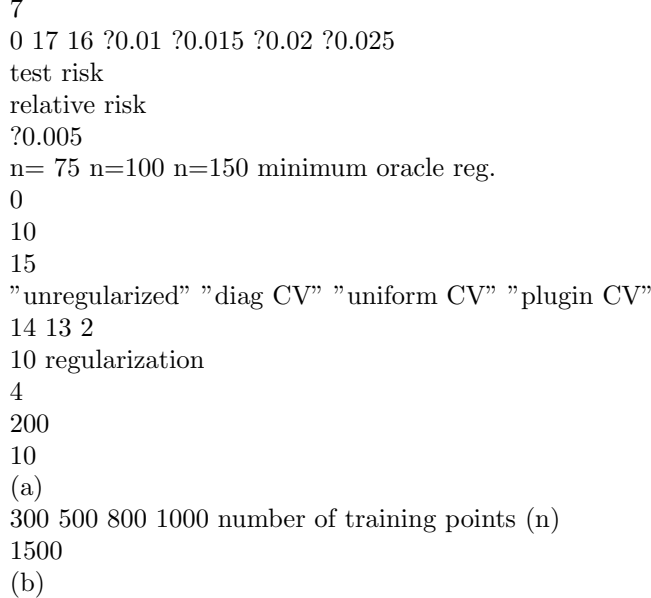


Figure 1: (a) Relative risk ($\text{Ln}(\gamma)$) of the hybrid generative/discriminative estimator for various γ ; the γ attaining the minimum of $\text{Ln}(\gamma)$ is close to the oracle γ^* (the vertical line). (b) On the MHC1 binding prediction task, test risk for the four multi-task estimators; P LUGIN CV (estimating all pairwise task affinities using P LUGIN and cross-validating the strength) works best.

30 independent train/test splits. Multi-task regularization actually performs worse than independent learning (D IAG CV) if we assume all tasks are equally related (U NIFORM CV). By learning the full matrix of task affinities (P LUGIN CV), we obtain the best results. Note that setting the $O(K^2)$ entries of γ via cross-validation is not computationally feasible, though other approaches are possible [13].

4 Related work and discussion

The subject of choosing regularization parameters has received much attention. Much of the learning theory literature focuses on risk bounds, which approximate the expected risk ($L(\gamma^*)$) with upper bounds. Our analysis provides a different type of approximation—one that is exact in the first few terms of the expansion. Though we cannot make a precise statement about the risk for any given n , exact control over the first few terms offers other advantages, e.g., the ability to compare estimators. To elaborate further, risk bounds are generally based on the complexity of the hypothesis class, whereas our analysis is based on the variance of the estimator. Vanilla uniform convergence

bounds yield worst-case analyses, whereas our asymptotic analysis is tailored to a particular problem (problem and algorithm (estimator)). Localization techniques [5], regret analyses [9], and stability-based bounds [8] all allow for some degree of problem- and algorithm-dependence. As bounds, however, they necessarily have some looseness, whereas our analysis provides exact constants, at least the ones associated with the lowest-order terms. Asymptotics has a rich tradition in statistics. In fact, our methodology of performing a Taylor expansion of the risk is reminiscent of AIC [1]. However, our aim is different: AIC is intended for model selection, whereas we are interested in optimizing regularization parameters. The Stein unbiased risk estimate (SURE) is another method of estimating the expected risk for linear models [21], with generalizations to non-linear models [11]. In practice, cross-validation procedures [10] are quite effective. However, they are only feasible when the number of hyperparameters is very small, whereas our approach can optimize many hyperparameters. Section 3.4 showed that combining the two approaches can be effective. To conclude, we have developed a general asymptotic framework for analyzing regularization, along with an efficient procedure for choosing regularization parameters. Although we are so far restricted to parametric problems with smooth losses and regularizers, we think that these tools provide a complementary perspective on analyzing learning algorithms to that of risk bounds, deepening our understanding of regularization. 8

2 References

- [1] H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19:716–723, 1974.
- [2] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *Advances in Neural Information Processing Systems (NIPS)*, pages 41–48, 2007.
- [3] B. Bakker and T. Heskes. Task clustering and gating for Bayesian multitask learning. *Journal of Machine Learning Research*, 4:83–99, 2003.
- [4] M. S. Bartlett. Approximate confidence intervals. II. More than one unknown parameter. *Biometrika*, 40:306–317, 1953.
- [5] P. L. Bartlett, O. Bousquet, and S. Mendelson. Local Rademacher complexities. *Annals of Statistics*, 33(4):1497–1537, 2005.
- [6] G. Bouchard. Bias-variance tradeoff in hybrid generative-discriminative models. In *Sixth International Conference on Machine Learning and Applications (ICMLA)*, pages 124–129, 2007.
- [7] G. Bouchard and B. Triggs. The trade-off between generative and discriminative classifiers. In *International Conference on Computational Statistics*, pages 721–728, 2004.
- [8] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [9] N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- [10] P. Craven and G. Wahba. Smoothing noisy data with spline functions. estimating the correct degree of smoothing by the method of generalized cross-validation. *Numerische Mathematik*, 31(4):377–403, 1978.
- [11] Y. C. Eldar. Generalized SURE for exponential families: Applications to regularization. *IEEE Transactions on Signal Processing*, 57(2):471–481, 2009.
- [12] T. Evgeniou,

C. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6:615?637, 2005. [13] L. Jacob, F. Bach, and J. Vert. Clustered multi-task learning: A convex formulation. In *Advances in Neural Information Processing Systems (NIPS)*, pages 745?752, 2009. [14] W. James and C. Stein. Estimation with quadratic loss. In *Fourth Berkeley Symposium in Mathematics, Statistics, and Probability*, pages 361?380, 1961. [15] J. A. Lasserre, C. M. Bishop, and T. P. Minka. Principled hybrids of generative and discriminative models. In *Computer Vision and Pattern Recognition (CVPR)*, pages 87?94, 2006. [16] P. Liang, F. Bach, G. Bouchard, and M. I. Jordan. Asymptotically optimal regularization in smooth parametric models. Technical report, ArXiv, 2010. [17] P. Liang and M. I. Jordan. An asymptotic analysis of generative, discriminative, and pseudolikelihood estimators. In *International Conference on Machine Learning (ICML)*, 2008. [18] A. McCallum, C. Pal, G. Druck, and X. Wang. Multi-conditional learning: Generative/discriminative training for clustering and classification. In *Association for the Advancement of Artificial Intelligence (AAAI)*, 2006. [19] B. Peters, H. Bui, S. Frankild, M. Nielson, C. Lundegaard, E. Kostem, D. Basch, K. Lamberth, M. Harndahl, W. Fleri, S. S. Wilson, J. Sidney, O. Lund, S. Buus, and A. Sette. A community resource benchmarking predictions of peptide binding to MHC-I molecules. *PLoS Computational Biology*, 2, 2006. [20] R. Raina, Y. Shen, A. Ng, and A. McCallum. Classification with hybrid generative/discriminative models. In *Advances in Neural Information Processing Systems (NIPS)*, 2004. [21] C. M. Stein. Estimation of the mean of a multivariate normal distribution. *Annals of Statistics*, 9(6):1135?1151, 1981. [22] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.