

Greedy Model Averaging

Authored by:

Tong Zhang
Dong Dai

Abstract

This paper considers the problem of combining multiple models to achieve a prediction accuracy not much worse than that of the best single model for least squares regression. It is known that if the models are misspecified, model averaging is superior to model selection. Specifically, let n be the sample size, then the worst case regret of the former decays at the rate of $O(1/n)$ while the worst case regret of the latter decays at the rate of $O(1/\sqrt{n})$. In the literature, the most important and widely studied model averaging method that achieves the optimal $O(1/n)$ average regret is the exponential weighted model averaging (EWMA) algorithm. However this method suffers from several limitations. The purpose of this paper is to present a new greedy model averaging procedure that improves EWMA. We prove strong theoretical guarantees for the new procedure and illustrate our theoretical results with empirical examples.

1 Paper Body

This paper considers the model combination problem, where the goal is to combine multiple models in order to achieve improved accuracy. This problem is important for practical applications because it is often the case that single learning models do not perform as well as their combinations. In practice, model combination is often achieved through the so-called "stacking" procedure, where multiple models $\{f_1(x), \dots, f_M(x)\}$ are first learned based on a shared "training dataset". Then these models are combined on a separate "validation dataset". This paper is motivated by this scenario. In particular, we assume that M models $\{f_1(x), \dots, f_M(x)\}$ are given a priori (e.g., we may regard them as being obtained with a separate training set), and we are provided with n labeled data points (validation data) $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ to combine these models. For simplicity and clarity, our analysis focuses on least squares regression in fixed design although similar analysis can be extended to random design and to other loss functions. In this setting, for notation convenience, we can represent the k -th model on the validation data as a vector $f_k = [f_k(X_1), \dots, f_k(X_n)]^T \in \mathbb{R}^n$, and we let the observation

vector $y = [Y_1, \dots, Y_n]^T \in \mathbb{R}^n$. Let $g = \mathbb{E}y$ be the mean. Our goal (in the fixed design or denoising setting) is to estimate the mean vector g from y using the M existing models $F = \{f_1, \dots, f_M\}$. Here, we can write $y = g + \epsilon$, where we assume that ϵ are iid Gaussian noise: $\epsilon \sim N(0, \sigma^2 I_n)$ for simplicity. This iid Gaussian assumption isn't critical, and the results remain the same for independent sub-Gaussian noise. We assume that the models may be mis-specified. That is, let k^* be the best single model defined as:

$$k^* = \arg \min_{k \in [M]} \|f_k - g\|_2^2,$$

$$(1)$$

then $\|f_{k^*} - g\|_2^2$. We are interested in an estimator \hat{f} of g that achieves a small regret

$$\frac{1}{n} \sum_{t=1}^n \|f_t - g\|_2^2$$

$$\frac{1}{n} \sum_{t=1}^n \|f_t - g\|_2^2$$

$R(\hat{f}) = \frac{1}{n} \sum_{t=1}^n \|f_t - g\|_2^2 - \frac{1}{n} \sum_{t=1}^n \|f_{k^*} - g\|_2^2$. This paper considers a special class of model combination methods which we refer to as model averaging, with combined estimators of the form $\hat{f} =$

$$\sum_{k=1}^M w_k f_k,$$

$$w_k \geq 0, \sum_{k=1}^M w_k = 1,$$

$$k=1$$

where $w_k \geq 0$ and $\sum_{k=1}^M w_k = 1$. A standard method for model averaging is model selection, where we choose the model k with the smallest least squares error: $\hat{f} = f_{k^*}$;

$$k^* = \arg \min_{k \in [M]} \|f_k - g\|_2^2.$$

However, it is well known that $w_{k^*} = 0$ when $k \neq k^*$. This corresponds to the choice of $w_{k^*} = 1$ and $w_k = 0$ for $k \neq k^*$. The worst case regret this procedure can achieve is $R(\hat{f}) = O(\ln M/n)$ [1]. Another standard model averaging method is the Exponential Weighted Model Averaging (EWMA) estimator defined as $\hat{f} = \sum_{k=1}^M w_k f_k$, where $w_k = \frac{e^{-\eta \sum_{j=1}^n \|f_k - g\|_2^2}}{\sum_{j=1}^M e^{-\eta \sum_{j=1}^n \|f_j - g\|_2^2}}$ with a tuned parameter $\eta > 0$. The extra parameters $\{q_j\}_{j=1, \dots, M}$ are priors that P impose bias favoring some models over some other models. Here we assume that $q_j \geq 0$ and $\sum_{j=1}^M q_j = 1$. In this setting, the most common prior choice is the flat prior $q_j = 1/M$. It should be pointed out that a progressive variant of (2), which returns the average of $n + 1$ EWMA estimators with $S_i = \{(X_1, Y_1), \dots, (X_i, Y_i)\}$ for $i = 0, 1, \dots, n$, was often analyzed in the earlier literature [2, 9, 5, 1]. Nevertheless, the non progressive version presented in (2) is clearly a more natural estimator, and this is the form that has been studied in more recent work [3, 6, 8]. Our current paper does not differentiate these two versions of EWMA because they have similar theoretical properties. In particular, our experiments only compare to the non-progressive version (2) that performs better in practice. It is known that exponential model averaging leads to an average regret of $O(\ln M/n)$ which achieves the optimal rate; however it was pointed out in [1] that the rate does not hold with large probability. Specifically, EWMA only leads to a sub-optimal deviation bound of $O(\ln M/n)$ with large probability. To remedy this sub-optimality, an empirical

star algorithm (which we will refer to as STAR from now on) was then proposed in [1]; it was shown that the algorithm gives $O(\ln M/n)$ deviation bound with large probability under the flat prior $q_i = 1/M$. One major issue of the STAR algorithm is that its average performance is often inferior to EWMA, as we can see from our empirical examples. Therefore although theoretically interesting, it is not an algorithm that can be regarded as a replacement of EWMA for practical purposes. Partly for this reason, a more recent study [7] re-examined the problem of improving EWMA, where different estimators were proposed in order to achieve optimal deviation for model averaging. However, the proposed algorithms are rather complex and difficult to implement. The purpose of this paper is to present a simple greedy model averaging (GMA) algorithm that gives the optimal $O(\ln M/n)$ deviation bound with large probability, and it can be applied with arbitrary prior q_i . Moreover, unlike STAR which has average performance inferior to EWMA, the average performance of GMA algorithm is generally superior to EWMA as we shall illustrate with examples. It also has some other advantages which we will discuss in more details later in the paper.

2

Greedy Model Averaging

This paper studies a new model averaging procedure presented in Algorithm 1. The procedure has L stages, and each time adds an additional model $f_k(\cdot)$ into the ensemble. It is based on a simple, but 2

important modification of a classical sequential greedy approximation procedure in the literature [4], which corresponds to setting $\eta(\cdot) = 0$, $\lambda = 0$ in Algorithm 1 with $\eta(\cdot)$ optimized over $[0, 1]$. The (2) STAR algorithm corresponds to the stage-2 estimator \hat{f}^* with the above mentioned classical greedy procedure of [4]. However, in order to prove the desired deviation bound, our analysis critically

2

(η 1)

η f which isn't present in the classical procedure (that depends on the extra term $\eta(\cdot) f_j$)

2

is, our proof does not apply to the procedure of [4]). As we will see in Section 4, this extra term does have a positive impact under suitable conditions that correspond to Theorem 1 and Theorem 2 below, and thus this term is not only for theoretical interest, but also it leads to practical benefits under the right conditions. Another difference between GMA and the greedy algorithm in [4] is that our procedure allows the use of non-flat priors through the extra penalty term $\eta(\cdot) \ln(1/q_j)$. This generality can be useful for some applications. Moreover, it is useful to notice that if we choose the flat prior $q_j = 1/M$, then the term $\eta(\cdot) \ln(1/q_j)$ is identical for all models, and thus this term can be removed from the optimization. In this case, the proposed method has the advantage of being parameter free (with the default choice of $\eta = 0.5$). This advantage is also shared by the STAR algorithm. : noisy observation y and static models f_1, \dots, f_M (\cdot) output : averaged model \hat{f} parameters: prior $\{q_j\}_{j=1,\dots,M}$ and regularization parameters η and λ input

the term $\sum_{j=1}^J \ln(1/q_j)$ from the estimators. This result also improves the recent work of [7] in that the resulting bound is scale free while the algorithm itself is significantly simpler. One disadvantage of this stage-2 estimator (and similarly the STAR estimator of [1]) is that its average performance is generally inferior to that of EWMA, mainly due to the relatively large constant in Theorem 1 (the same issue holds for the STAR algorithm). For this reason, the stage-2 estimator is not a practical replacement of EWMA. This is the main reason why it is necessary to run GMA for $L \geq 2$ stages, which leads to reduced constants (see Theorem 2) below. Our empirical experiments show that in order to compete with EWMA for average performance, it is important to take $L \geq 2$. However a relatively small L (as small as $L = 5$) is often sufficient, and in such case the resulting estimator is still quite sparse. \square

Theorem 1 Given $q_j \geq 0$ such that

$q_j = 1$. If $\sum_{j=1}^J q_j \geq 2$, then with probability $1 - \delta$ we have

$\sum_{j=1}^J$

$R(f^*)$

(2)

)

$\leq 3 \sum_{j=1}^J \ln(1/q_j) + \ln(1/\delta) \leq n \sum_{j=1}^J q_j$

(2) While the stage-2 estimator f^* achieves the optimal rate, running GMA for more more stages can further improve the performance. The following theorem shows that similar bounds can be obtained for GMA at stages larger than 2. However, the constant before $\sum_{j=1}^J \ln(1/q_j)$ approaches 8 as $\delta \rightarrow 0$ (with default $\delta = 0.5$), which is smaller than the constant of Theorem 1 which is about 30. This implies potential improvement when we run more stages, and this improvement is confirmed in our empirical study. In fact, with relatively large δ , the GMA method not only has the theoretical advantage of achieving smaller regret in deviation (that is, the regret bound holds with large probability) but also achieves better average performance in practice. \square

Theorem 2 Given $q_j \geq 0$ such that

$q_j = 1$. If $\sum_{j=1}^J q_j \geq 2$ and let $0 \leq \delta \leq 1$ in Algorithm 1,

$\sum_{j=1}^J$

then with probability $1 - \delta$ we have

$\sum_{j=1}^J \ln(1/q_j) \leq \ln(1/\delta) + 30 \sum_{j=1}^J \ln(1/q_j) \leq n \sum_{j=1}^J q_j$

Another important advantage of running GMA for $L \geq 2$ stages is that the resulting estimator not only competes with the best single estimator f^* , but also competes with the best estimator in the convex hull of $\text{cov}(F)$ (with the parameter δ appropriately tuned). Note that the latter can be significantly better than the former. Define the convex hull of F as $\text{CH}(F) = \{ \sum_{j=1}^J w_j f_j : w_j \geq 0; \sum_{j=1}^J w_j = 1 \}$.

(*) The following theorem shows that as $\delta \rightarrow 0$, the prediction error of f^* is no more than $O(1/n)$ worse than that of the optimal $f^* \in \text{CH}(F)$ when we choose a sufficiently small $\delta = O(1/n)$ in Algorithm 1. Note that in this case, it is beneficial to use a parameter δ smaller than the default choice of $\delta = 0.5$. This phenomenon is also confirmed by our experiments. \square

$q_j \geq 0$ such that $\sum_j q_j = 1$. Consider any $\{w_j : j = 1, \dots, M\}$ such that $\sum_{j=1}^M w_j = 1$ and $w_j \geq 0$, and let $f = \sum_{j=1}^M w_j f_j$. If $\epsilon \leq 40\epsilon$ and let $0 \leq \epsilon \leq 1$ in Algorithm 1, then with probability $1 - 2\epsilon$, when $\epsilon \leq \epsilon$:

$$\begin{aligned} & \frac{1}{n} \sum_{t=1}^n \sum_{k=1}^K \sum_{j=1}^M w_j \sum_{i=1}^n (f_i - f_{j,t})^2 \\ & \leq \frac{1}{n} \sum_{t=1}^n \sum_{k=1}^K \sum_{j=1}^M w_j \sum_{i=1}^n (f_i - f_j)^2 + O\left(\frac{1}{n}\right) \\ & \leq \frac{1}{n} \sum_{t=1}^n \sum_{k=1}^K \sum_{j=1}^M w_j \sum_{i=1}^n (f_i - f_j)^2 + O\left(\frac{1}{n}\right) \\ & \leq \frac{1}{n} \sum_{t=1}^n \sum_{k=1}^K \sum_{j=1}^M w_j \sum_{i=1}^n (f_i - f_j)^2 + O\left(\frac{1}{n}\right) \end{aligned}$$

Experiments

The point of these experiments is to show that the consequences of our theoretical analysis can be observed in practice, which support the main conclusions we reach. For this purpose, we consider the model $g = Xw + 0.5g$, where $X = (f_1, \dots, f_M)$ is an $n \times M$ matrix with independent standard Gaussian entries, and $g \sim N(0, I_n)$ implies that the model is mis-specified. The noise vector is $\epsilon \sim N(0, \frac{1}{2} I_n)$, independently generated of X . The coefficient vector $w = (w_1, \dots, w_M)$ is given by $w_i = \frac{1}{M} \sum_{j=1}^M u_j$ for $i = 1, \dots, M$, where u_1, \dots, u_M are independent standard uniform random variables for some fixed s . The performance of an estimator \hat{f} measured here is the mean squared error (MSE) defined as

$$\frac{1}{n} \sum_{t=1}^n \sum_{k=1}^K \sum_{j=1}^M w_j \sum_{i=1}^n (f_i - \hat{f}_{j,t})^2$$

We run the Greedy Model Averaging (GMA) algorithm for L stages up to $L = 40$. The EWMA parameter is tuned via 10-fold cross-validation. Moreover, we also listed the performance of EWMA with projection, which is the method that runs EWMA, but with each model f_k replaced by model $\hat{f}_k = \sum_{j=1}^M \hat{w}_{j,k} f_j$ where $\hat{w}_k = \arg \min_{w \geq 0, \sum w_j = 1} \sum_{i=1}^n (f_i - \sum_{j=1}^M w_j f_j)^2$. That is, \hat{f}_k is the best linear scaling of f_k to predict y . Note that this is a special case of the class of methods studied in [6] (which considers more general projections) that leads to non progressive regret bounds, and this is the method of significant current interests [3, 8]. However, at least for the scenario considered in our paper, the projected EWMA method never improves performance in our experiments. Finally, for reference purpose, we also report the MSE of the best single model (BSM) \hat{f}_k , where k^* is given by (1). The model \hat{f}_k is clearly not a valid estimator because it depends on the unobserved g ; however its performance is informative, and thus included in the tables. For simplicity, all algorithms use flat prior $q_k = 1/M$.

4

Illustration of Theorem 1 and Theorem 2

The first set of experiments are performed with the parameters $n = 50$, $M = 200$, $s = 1$ and $\epsilon = 2$. Five hundred replications are run, and the MSE performance of different algorithms are reported in Table 1 using the $\bar{\epsilon}$ mean ϵ .

standard deviation? format. Note that with $s = 1$, the target is $g = f_1 + 0.5?g$. Since f_1 and $?g$ are random Gaussian vectors, the best single model is likely f_1 . The noise $? = 2$ is relatively large. This is thus the situation that model averaging does not achieve as good a performance as that of the best single model. This corresponds to the scenario considered in Theorem 1 and Theorem 2. The results indicate that for GMA, from $L = 1$ (corresponding to model selection) to $L = 2$ (stage-2 model averaging of Theorem 1), there is significant reduction of error. The performance of GMA with $L = 2$ is comparable to that of the STAR algorithm. This isn't surprising, because STAR can be regarded as the stage-2 estimator based on the more classical greedy algorithm of [4]. We also observe that the error keeps decreasing (but at a slower pace) when $L \geq 2$, which is consistent with Theorem 2. It means that in order to achieve good performance, it is necessary to use more stages than $L = 2$ (although this doesn't change the $O(1/n)$ rate for regret, it can significantly reduce constant). It becomes better than EWMA when L is as small as 5, which still gives a relatively sparse averaged model. EWMA with projection does not perform as well as the standard EWMA method in this setting. Moreover, we note that in this scenario, the standard choice of $? = 0.5$ in Theorem 2 is superior to choosing smaller $? = 0.1$ or $? = 0.001$. This again is consistent with Theorem 2, which shows that the new term we added into the greedy algorithm is indeed useful in this scenario.

5

Illustration of Theorem 3

The second set of experiments are performed with the parameters $n = 50$, $M = 200$, $s = 10$ and $? = 0.5$. Five hundred replications are run, and the MSE performance of different algorithms are reported in Table 2 using the 'mean ? standard deviation' format. 5

Table 1: MSE of different algorithms: best single model is superior to averaged models STAR EWMA EWMA (with projection) BSM 0.663 ? 0.4 0.645 ? 0.5 0.744 ? 0.5 0.252 ? 0.05 GMA ? = 0.5 ? = 0.1 ? = 0.01

L=1 0.735 ? 0.74 0.735 ? 0.74 0.735 ? 0.74

L=2 0.689 ? 0.4 0.689 ? 0.4 0.689 ? 0.4

L=5 0.58 ? 0.39 0.645 ? 0.31 0.663 ? 0.3

L = 20 0.566 ? 0.37 0.623 ? 0.29 0.638 ? 0.28

L = 40 0.567 ? 0.38 0.622 ? 0.29 0.639 ? 0.28

Note that with $s = 10$, the target is $g = f? + 0.5?g$ for some $f? \in \text{cov}(F)$. The noise $? = 0.5$ is relatively small, which makes it beneficial ? to compete with the best model $f?$ in the convex hull even though GMA has a larger regret of $O(1/n)$ when competing with $f?$. This is thus the situation considered in Theorem 3, which means that model averaging can achieve better performance than that of the best single model. The results again show that for GMA, from $L = 1$ (corresponding to model selection) to $L = 2$ (stage-2 model averaging of Theorem 1), there is significant reduction of error. The performance of GMA with $L = 2$ is again comparable to that of the STAR algorithm. Again we observe that even with the standard choice of $? = 0.5$, the error keeps decreasing (but at a slower pace) when $L \geq 2$, which is consistent with Theorem 2. It becomes

better than EWMA when L is as small as 5, which still gives a relatively sparse averaged model. EWMA with projection again does not perform as well as the standard EWMA method in this setting. Moreover, we note that in this scenario, the standard choice of $\gamma = 0.5$ in Theorem 2 is inferior to choosing smaller parameter values of $\gamma = 0.1$ or $\gamma = 0.001$. This is consistent with Theorem 3, where it is beneficial to use a smaller value for γ in order to compete with the best model in the convex hull.

Table 2: MSE of different algorithms: best single model is inferior to averaged model STAR EWMA EWMA (with projection) BSM 0.443 0.08 0.316 0.087 0.364 0.078 0.736 0.083 GMA $\gamma = 0.5$ $\gamma = 0.1$ $\gamma = 0.01$

6

$L=1$ 0.809 0.12 0.809 0.12 0.809 0.12

$L=2$ 0.456 0.081 0.456 0.081 0.456 0.081

$L=5$ 0.305 0.062 0.269 0.056 0.268 0.053

$L = 20$ 0.266 0.057 0.214 0.046 0.211 0.045

$L = 40$ 0.265 0.057 0.211 0.045 0.207 0.045

Conclusion

This paper presents a new model averaging scheme which we call greedy model averaging (GMA). It is shown that the new method can achieve regret bound of $O(\ln M/n)$ with large probability when competing with the single best model. Moreover, it can also compete with the best combined model in convex hull. Both our theory and experimental results suggest that the proposed GMA algorithm is superior to the standard EWMA procedure. Due to the simplicity of our proposal, GMA may be regarded as a valid alternative to the more widely studied EWMA procedure both for practical applications and for theoretical purposes. Finally we shall point out that while this work only considers static model averaging where the models F are finite, similar results can be obtained for affine estimators or infinite models considered in recent work [3, 6, 8]. Such extension will be left to the extended report.

A

Proof Sketches

We only include proof sketches, and leave the details to the supplemental material that accompanies the submission. First we need the following standard Gaussian tail bounds. The proofs can be found in the supplemental material. 6

Proposition 1 Let $f_j \in \mathbb{R}^n$ be a set of fixed vectors ($j = 1, \dots, M$), and assume that $q_j \geq 0$ with $\sum_{j=1}^M q_j = 1$. Let k^* be a fixed integer between 1 and M . Define event E_1 as

$E_1 = \{j : (f_j - f_{k^*})^\top (f_j - f_{k^*}) \leq k^2 \ln(1/(q_j))\}$ and define event E_2 as

$E_2 = \{j, k : (f_j - f_k)^\top (f_j - f_k) \leq k^2 \ln(1/(q_j q_k))\}$, then $P(E_1) \geq 1 - 2^{-k^2}$ and $P(E_2) \geq 1 - 2^{-k^2}$. A.1

Proof Sketch of Theorem 1

More detailed proof can be found in the supplemental material. Note that with probability $1 - 2^{-k^2}$, both event E_1 and event E_2 of Proposition 1 hold. Moreover we have

2

$$\begin{aligned}
& 2 \\
& (2) \quad (1) \\
& ? \\
& f ? g = ?(2) f ? + (1 ? ?(2)) f k ?(2) ? g 2 2 \\
& 2 \\
& (2) ? (1) (2) ? ? f + (1 ? ?) f k ? ? g + 2(1 ? ?(2)) ? \leq (f k ?(2) ? f k ?) 2 \\
& 2 (1) \\
& 2 (1) \\
& (2) ? ? + ? \\
& f ? f k ? ? f ? f k ?(2) + ?c(2) (\ln(1/qk ?) ? \ln(1/qk ?(2))). 2 \\
& 2 \\
& (2) \\
& ?(2)
\end{aligned}$$

In the above derivation, the inequality is equivalent to $Q(k) \leq Q(2)(k)$, which is a simple fact of the definition of $k^{(\cdot)}$ in the algorithm. Also we can rewrite the fact that $Q(1)(k^{(1)}) \leq Q(1)(k)$ as

$$\begin{aligned}
& (1) \\
& 2 \\
& 2 \\
& ? \\
& f ? g ? f k ? ? g 2 ? 2 ? \leq (f k ?(1) ? f k ?) + ?c(1) \ln(qk ?(1) / qk ?). 2 \\
& \text{By combining the above two inequalities, we obtain}
\end{aligned}$$

$$\begin{aligned}
& (2) \\
& 2 \text{ i h} \\
& 2 \\
& ? \\
& f ? g ? f k ? ? g 2 ? ?(2) 2 ? \leq (f k ?(1) ? f k ?) + ?c(1) \ln(qk ?(1) / qk ?) 2 \\
& \text{h i} \\
& 2 + 2(1 ? ?(2)) ? \leq (f k ?(2) ? f k ?) + ?(2) ? ?(2) (1 ? ?(2)) f k ?(1) ? f k ? \\
& 2 \\
& 2 ? ?(2) f k ?(1) ? f k ?(2) 2 + ?c(2) (\ln(1/qk ?) ? \ln(1/qk ?(2))). \text{ Since } ?(2) \\
& = 1/2, \text{ we obtain}
\end{aligned}$$

$$\begin{aligned}
& (2) \\
& 2 \\
& 2 \\
& ? \\
& f ? g ? f k ? ? g 2 2 \\
& 1 1 ? (?c(1) + ?c(2)) \ln(1/qk ?) ? ?c(1) \ln(1/qk ?(1)) ? ?c(2) \ln(1/qk ?(2)) \\
&) 2 2 \text{ s s} \\
& 1 1 1 \\
& + 2 f k ?(1) ? f k ? 2 ? 2 \ln + 2 ? f k ?(2) ? f k ?(1) 2 ? 2 \ln qk ?(1) ? 2 qk ?(1) \\
& qk ?(2) ? \\
& 2 2 + ?(2) ? 1/4 f k ?(1) ? f k ? 2 ? ?(2) f k ?(1) ? f k ?(2) 2 1 ? (?c(1) \\
& + ?c(2)) \ln(1/qk ?) + (2r_1 + 2r_2) \ln(1/?). 2 \text{ The first inequality above uses} \\
& \text{the tail probability bounds in the event E1 and E2 . We then use the algebraic}
\end{aligned}$$

Proof Sketch of Theorem 2

Now consider any $i \geq 3$. We have

2

(6)

2 i h ('?1)

?

$$+ (1 - \beta) \frac{1}{\beta} \left(\frac{1}{\beta} \right) f(k) + 2\beta \frac{1}{\beta} \left(\frac{1}{\beta} \right) f(k) + (1 - \beta) \frac{1}{\beta} \left(\frac{1}{\beta} \right) f(k)$$

f ? gg ? ?(') f? 2 2

2

2 ('?1)

$$(\text{?1}) \text{ (')} \text{ ? f ?(')} . +?c (\ln(1/qk) \text{ ? } \ln(1/q?('))) + ?(') \text{ f? ? f ? f? ?}$$

 (\cdot)

?

 $k?$
$$k$$
$$\mathbf{k}$$

2

2

$$\mathcal{L}(\mathbf{y})$$
$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

2

2

?

$$f^? \ g^? \ f \ k^? \ ? \ g \ 2 \ 2$$

2

2

$$(\cdot) \quad ? \quad (\cdot ? 1)$$

[illegible]

f 2

2

hi ('?1) ('?1) 2

? f k? ??(') f k?(') ? f?

$$+ \frac{1}{2} \left(\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} \right) = \frac{1}{2} \Delta f$$

2

(?1) 2

$$\frac{1}{2} \left(\frac{1}{qk} + \frac{1}{qk'} \right) \ln \left(\frac{1}{qk} \right) + \frac{1}{2} \left(\frac{1}{qk} + \frac{1}{qk'} \right) \ln \left(\frac{1}{qk'} \right) + \frac{1}{2} \left(\frac{1}{qk} + \frac{1}{qk'} \right) \ln \left(\frac{1}{qk} \right) + \frac{1}{2} \left(\frac{1}{qk} + \frac{1}{qk'} \right) \ln \left(\frac{1}{qk'} \right)$$

$$?(') ?(') (1 ? ?(')) ? ?(') 2 1$$

$$f^{(i)}(k) = f(k) + f(k^2) + f(k^3) + \dots + f(k^i) = 2 + 1 + \dots + 1 = i$$

